



**HAL**  
open science

# Survival probability and size of lineages in antibody affinity maturation

Marco Molari, Rémi Monasson, Simona Cocco

► **To cite this version:**

Marco Molari, Rémi Monasson, Simona Cocco. Survival probability and size of lineages in antibody affinity maturation. *Physical Review E*, 2021, 103 (5), pp.052413. 10.1103/PhysRevE.103.052413 . hal-02974976v2

**HAL Id: hal-02974976**

**<https://hal.science/hal-02974976v2>**

Submitted on 10 May 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Survival probability and size of lineages in antibody affinity maturation

Marco Molari, Rémi Monasson, and Simona Cocco\*  
*Laboratoire de Physique de l'École Normale Supérieure, ENS,  
PSL University, CNRS UMR8023, Sorbonne Université,  
Université de Paris, 24 rue Lhomond, 75005 Paris, France.*  
(Dated: May 10, 2021)

Affinity Maturation (AM) is the process through which the immune system is able to develop potent antibodies against new pathogens it encounters, and is at the base of the efficacy of vaccines. At its core AM is analogous to a Darwinian evolutionary process, where B-cells mutate and are selected on the base of their affinity for an Antigen (Ag), and Ag availability tunes the selective pressure. In cases when this selective pressure is high the number of B-cells might quickly decrease and the population might risk extinction in what is known as a *population bottleneck*. Here we study the probability for a B-cell lineage to survive this bottleneck scenario as a function of the progenitor affinity for the Ag. Using recursive relations and probability generating functions we derive expressions for the average extinction time and progeny size for lineages that go extinct. We then extend our results to the full population, both in the absence and presence of competition for T-cell help, and quantify the population survival probability as a function of Ag concentration and initial population size. Our study suggests the population bottleneck phenomenology might represent a limit case in the space of biologically plausible maturation scenarios, whose characterization could help guide the process of vaccine development.

## I. INTRODUCTION

Affinity Maturation (AM) is a biological process through which our immune system generates potent Antibodies (Ab) against newly-encountered pathogens. AM is also at the base of the efficacy of vaccination, in which this process is artificially elicited through the administration of a dose of Antigen (Ag). The biological mechanisms that govern AM are many and complex, and are the object of many excellent reviews [1–8]. Simply speaking, AM works by subjecting a population of B-lymphocytes (or B-cells) to iterative cycles of mutations and selection for Ag binding, which generate a Darwinian evolutionary process that progressively increases their affinity for the Ag. This process is schematically depicted in fig. 1. Maturation takes place in Germinal Centers (GCs), microanatomical structures that appear inside of secondary lymphoid organs. They are divided in two areas: the GC Dark Zone (DZ) in which cells divide and mutate,<sup>1</sup> and the Light Zone (LZ) in which they undergo selection. Cells iteratively migrate between these two compartments. Selection in the LZ is completed in two steps. In the first step cells try to bind the Ag, exposed on the surface of Follicular Dendritic Cells (FDCs). In the second step they compete to receive a survival signal from T-follicular helper (Tfh) cells, in the absence of which they undergo apoptosis. Tfh cells are able to probe the amount of Ag captured by B-cells, preferentially delivering this survival signal to the cells that were most successful in capturing Ag. Cells that receive the signal either

migrate back to the DZ for additional rounds of mutation and selection, or they can differentiate in Plasma or Memory Cells (PCs/MCs). The former are responsible for the production of Abs to fight the infection, while the latter confer long-lasting protection by remaining quiescent until the same Ag is encountered again, in which case they reactivate and produce Abs or enter GCs for further maturation.

In spite of the many recent experimental advancements in the study of AM, several open questions still remain to be answered, which have important implications in vaccine design. For example understanding the role of Ag availability in controlling maturation might lead to optimization of Ag dosage in vaccines [12–14], or understanding the effect of B-cell precursor frequency and affinity might help improving immunogen design [15, 16]. Given the complexity of this process, computational models represent an invaluable tool to guide our understanding of AM [17, 18]. In this paper we introduce a stochastic model of AM to study the survival probability of B-cell lineages in GCs. Experimental analysis of vaccine-responsive lineages shows signatures of selection in their reconstructed phylogenies [19]. This selection pressure, which is partially controlled by Ag availability [13], is important to push lineages towards maturation, but at the same time an excessive pressure might be deleterious. Indeed, several maturation models present a phenomenology termed *population bottleneck* [20–22], in which strong selection pressure causes a decrease in GC population size, potentially leading to extinction. As a consequence of this trade-off optimal maturation is achieved at intermediate levels of selection pressure. While the bottleneck phenomenology has been studied through numerical simulations so far, we think that it is not incompatible with experimental evidence [23], and might represent a limiting regime for AM as

---

\* Email: simona.cocco@phys.ens.fr

<sup>1</sup> In the DZ B-cells express high levels of *Activation-Induced cytidine Deaminase*, an enzyme that increases the natural rate of DNA mutations up to  $10^{-3}$  per base-pair per generation [9–11].

later discussed in section V A.

The study of the survival probability of a population through the maturation bottleneck we report below is both analytical and numerical. We start by considering the dependence of a lineage survival probability on the progenitor affinity. Through the use of recursive relations and probability generating functions we are able to evaluate this probability, and also quantify extinction time and progeny size for lineages that go extinct. We then extend our approach to analyze the extinction probability for the full B-cell population, and its dependence on Ag concentration and initial population size. Last of all, we discuss the biological relevance of the bottleneck phenomenology, and how quantifying lineage survival probability might help vaccine design.

## II. MODEL FOR STOCHASTIC MATURATION

Our model for stochastic maturation is inspired by previous works [13, 21]. The model is simple enough to be analytically tractable, while retaining the main aspects of the bottleneck phenomenology.

### A. Steps in affinity maturation

We consider the evolution of a population of B-cells inside a GC. Through repeated cycles of mutation and

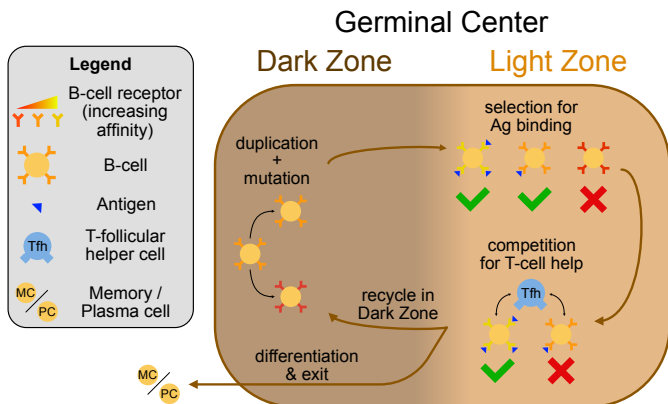


FIG. 1. schematic depiction of the germinal center reaction. Inside of a Germinal Center (GC) B-cells undergo iterative cycles of duplication, mutation and selection. Cell duplication and mutation occurs in the GC dark zone, while selection takes place in the light zone. Selection is divided in two steps: cells must first bind the antigen with their B-cell receptors, and then compete to receive a survival signal from T-follicular helper cells. Failure in any of these steps results in apoptosis. Successfully selected cells have some probability of differentiating into Memory or Plasma Cells (MCs/PCs), and exit the GC. Cells that do not differentiate recycle back in the dark zone to start a new maturation cycle.

selection the population increases its average affinity for the Ag over time. In our model each cell in the population is solely characterized by its affinity for the Ag, measured in terms of binding energy  $\epsilon$  and expressed in units of  $k_B T$ .

The simulation starts when the GC is mature (roughly 1 week after Ag injection [2]). The initial population is composed of  $N_i$  cells whose binding energy is independently extracted from a Gaussian distribution of naive responders, with mean  $\mu_i$  and standard deviation  $\sigma_i$ . Cells undergo iterative rounds of duplication, mutation and selection. These steps are schematized in fig. 2.

At the beginning of the round cells duplicate once in the GC DZ. Each daughter cell can then independently either:

- undergo an affinity-affecting mutation with probability  $p_{aa}$ , which causes its binding energy to change by some amount  $\Delta\epsilon$ . We assume that  $\Delta\epsilon$  is a random variable, extracted from a Gaussian distribution with mean  $\mu_M$  and standard deviation  $\sigma_M$  (see fig. 2 “mutation”);
- not mutate or develop silent mutations, with probability  $p_{sil}$ . In both cases its affinity is unchanged;
- be hit by a lethal mutation with probability  $p_{let}$ , in which case the cell is removed from the population.

As a result, the total distribution of the changes  $\Delta\epsilon$  is therefore summarized by the kernel

$$K(\Delta\epsilon) = \frac{p_{aa}}{\sqrt{2\pi\sigma_M^2}} \exp\left(-\frac{(\Delta\epsilon - \mu_M)^2}{2\sigma_M^2}\right) + p_{sil} \delta(\Delta\epsilon) \quad (1)$$

where  $\delta(\Delta\epsilon)$  is Dirac delta distribution. Notice that, due to lethal mutations, the integral of the kernel  $K$  is not normalized to unity but to  $p_{aa} + p_{sil} = 1 - p_{let}$ . Parameters are chosen such that only a small fraction of the mutations is beneficial, i.e. decreases the binding energy (cf. appendix A).

After duplication and mutation cells migrate to the LZ where they try to bind the Ag exposed on the surface of FDCs. Failure to do so results in cell death, and only cells that are able to bind the Ag with sufficient affinity survive this step of selection. Similarly to [13, 21] we consider the survival probability for a cell with binding energy  $\epsilon$  to be given by the following Langmuir isotherm:

$$P_{Ag}(\epsilon) = \frac{C e^{-\epsilon}}{C e^{-\epsilon} + e^{-\epsilon_{Ag}}}, \quad (2)$$

where  $\epsilon_{Ag}$  is a threshold binding energy and  $C$  represents the dimensionless concentration of Ag available for cells to bind. This concentration controls the strength of selection, making successful binding more likely when more Ag is available to bind. In practice it acts by imparting a shift of magnitude  $\log C$  to the energy threshold. The functional dependence of the selection probability on  $\epsilon$  and  $C$  is displayed in fig. 2 (“selection for Ag binding”).

In a second selection step cells compete to receive a survival signal from T-follicular helper cells, with the signal being preferentially delivered to cells that bind more Ag. The survival probability for a cell with binding energy  $\epsilon$  is:

$$P_T(\epsilon, \bar{\epsilon}) = \frac{C e^{-\epsilon}}{C e^{-\epsilon} + e^{-\bar{\epsilon}}}, \quad \text{with } e^{-\bar{\epsilon}} = \langle e^{-\epsilon} \rangle_{\text{pop}} \quad (3)$$

Where the term  $\langle e^{-\epsilon} \rangle_{\text{pop}}$  represents the average of this quantity over the population and encodes for competition (see fig. 2 “competitive selection for T-cell help”). The surviving cells can then differentiate into plasma or memory cells with total probability  $p_{\text{diff}}$ . We do not keep track of these differentiated cells in the simulation.

After this step if the population size exceeds the maximum carrying capacity  $N_{\text{max}}$  cells are randomly removed until this threshold is met. The surviving cells start then the next round of evolution. The values of the model parameters are reported in table I, and discussed in appendix A.

### B. Population bottleneck and lineage survival

Similarly to other AM models [20, 21], for standard parameter values the population initially undergoes a bottleneck state. This is caused by the strong selection pressure initially imposed by Ag-binding selection, which later relaxes if the average population energy reaches values  $\langle \epsilon \rangle_{\text{pop}} \sim \epsilon_{\text{Ag}}$ . By controlling the selection pressure (cf. eqs. (2) and (3)) Ag concentration also impacts the population survival probability.

As an illustration we report in fig. 3 the average evolution of 1000 stochastic simulations for three different values of the concentration  $C$ . For all three values the population size initially decreases under the combined effect of the two selection steps (fig. 3A). This decrease lasts for few evolution rounds, and is accompanied by a quick increase in average affinity (fig. 3B). At this point surviving populations are composed of few high-affinity cells, on which the main acting selection force is competitive selection in eq. (3). If this selection pressure is not too strong then the population will later expand and mature. Through a mechanism analogous to the one studied in [13] Ag concentration then controls the maturation speed, as can be seen by comparing the speed of decrease in average binding energy of the population after the bottleneck in the cases  $C = 10$  and  $C = 3.5$  in fig. 3B.

The fraction of surviving simulations as a function of time is shown in fig. 3C. At low concentration ( $C = 1$ ) the population goes quickly extinct in all simulations. For such small values of Ag concentration competitive selection alone is sufficiently strong to impede population growth. Intermediate concentration values ( $C = 3.5$ ), on the contrary, are sufficient to sustain population growth. In this case extinction can nevertheless occur close to the bottleneck state, when population size gets transiently

small, see fig. 3A; if some cells are able to survive and pass this bottleneck, then the population again grows to full size and continues maturation. Last of all, at high concentration ( $C = 10$ ), the bottleneck pressure is not sufficient to significantly endanger population survival, and all simulations are able to overcome the low-population state without going extinct, but maturation proceeds more slowly. We will study in detail in the next section the dependence of the survival probability of the population of cells on Ag concentration.

Survival and future expansion is also strongly dependent on the initial distribution of affinities. This effect can be readily observed on lineages originated from a single ancestor, with energy  $\epsilon$ . In fig. 4 we display three examples of lineage evolution in the form of trees in which each node corresponds to a different cell. These lineages differ by the affinity of their progenitor at the root of the tree. The progeny of the lowest-affinity one (red,  $\epsilon_i \sim -0.3$ ) goes extinct in few evolution rounds. In the one with intermediate affinity (orange,  $\epsilon_i \sim -0.45$ ) only few individuals are able to survive the bottleneck. The high-affinity one (green,  $\epsilon_i \sim -1.3$ ) instead expands and eventually takes over the population. To quantify the population survival probability we will first investigate how the survival of single lineages depends on the progenitor affinity.

## III. PROBABILITY OF SURVIVAL AND DISTRIBUTION OF EXTINCTION TIMES

In this section we study the probability that a B-cell lineage descending from a single progenitor cell survives through a population bottleneck, in particular how the probability of survival depends on the affinity of the progenitor. We also determine the distribution of extinction times of the lineage. We then make use of these results to evaluate the survival probability for the full population.

### A. Case of one lineage

#### 1. Probability of survival

Let us consider a progenitor cell with binding energy  $\epsilon$ , present in the population at the beginning of evolution  $t = 0$ . At each evolution round this cell will divide and its offspring will have some probability of being removed from the population, either due to selection, differentiation or lethal mutations. Since we are interested in studying the bottleneck phase, in which the population is not at its maximum size, we can neglect the enforcement of the carrying capacity constraint. In fig. 5A we report an example of lineage evolution for a progenitor with binding energy  $\epsilon = 1$ . Color indicates the binding energy of each cell, according to the color-scale on top. In this example cells accumulate deleterious mutations until

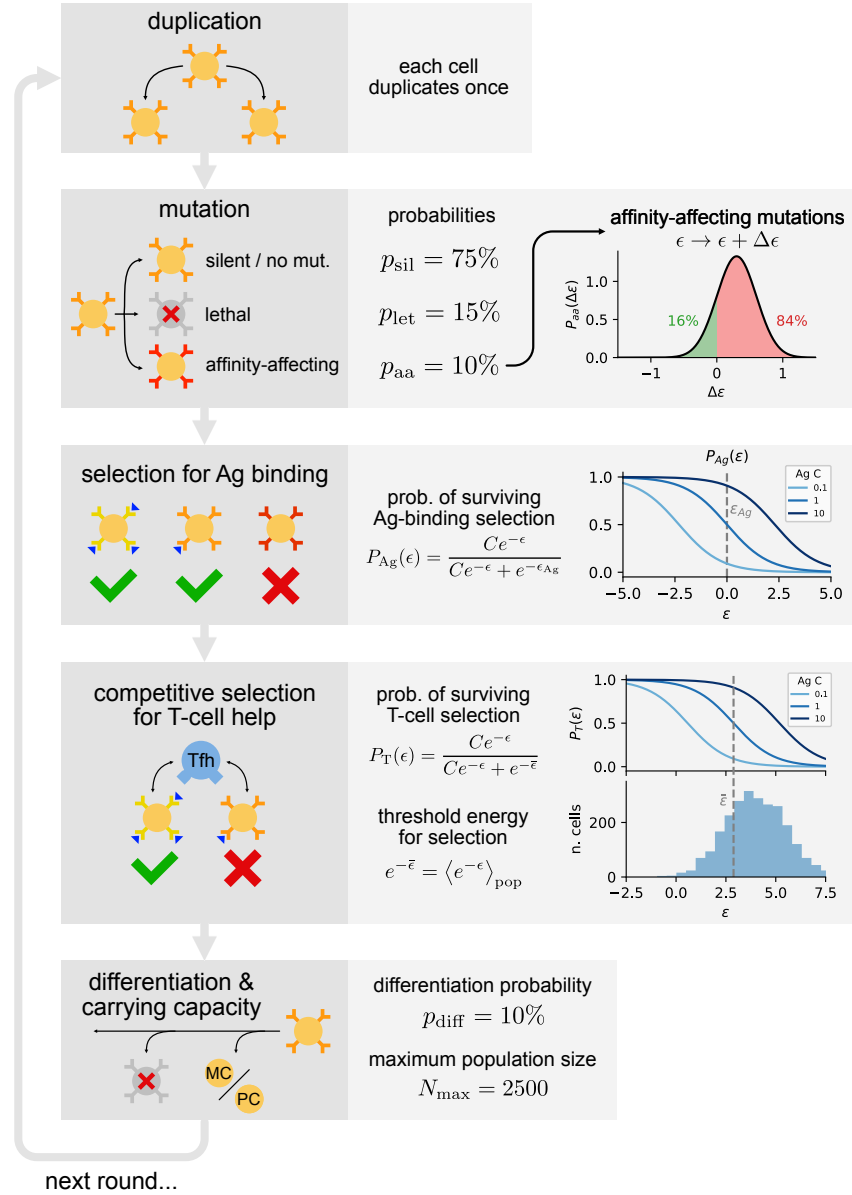


FIG. 2. schematic description of the processes that make up a simulated evolution round in our model. At the beginning of the round cells duplicate once. Each cell can then independently develop a mutation. Cells that undergo a lethal mutation ( $p_{let} = 15\%$ ) are removed from the population, while cells that develop an affinity-affecting mutation ( $p_{aa} = 10\%$ ) receive an additive change  $\Delta\epsilon$  to their binding energy, extracted from the displayed Gaussian distribution. Notice that most of the mutations have a deleterious effect on affinity. Cells undergo then selection for antigen binding and compete to receive T-cell help. For each selection step we display the functional behavior of the probability of surviving as a function of the progenitor affinity  $\epsilon$  and antigen concentration  $C$ . In selection for T-cell help competition is obtained by making the threshold energy  $\bar{\epsilon}$  depend on the current affinity distribution of the population, as displayed. Cells that are able to survive selection have a probability  $p_{diff} = 10\%$  of differentiating into memory or plasma cells, and exiting the cycle. Last of all, if the population size exceeds the threshold value  $N_{max} = 2500$  cells in excess are randomly removed. The remaining cells will then begin a new evolution round.

the lineage eventually goes extinct after  $t = 26$  evolution rounds.

We are interested in computing the probability  $d_t(\epsilon)$  that *all of the offspring* of a progenitor with binding energy  $\epsilon$  will be extinct by evolution round  $t$ , see fig. 5B.

The expression for  $t = 1$  can easily be written as the probability that both daughter cells generated during the duplication phase will be removed by the end of the round. As stated above, this can occur either by lethal mutation, by failing selection or by differentiation. For each

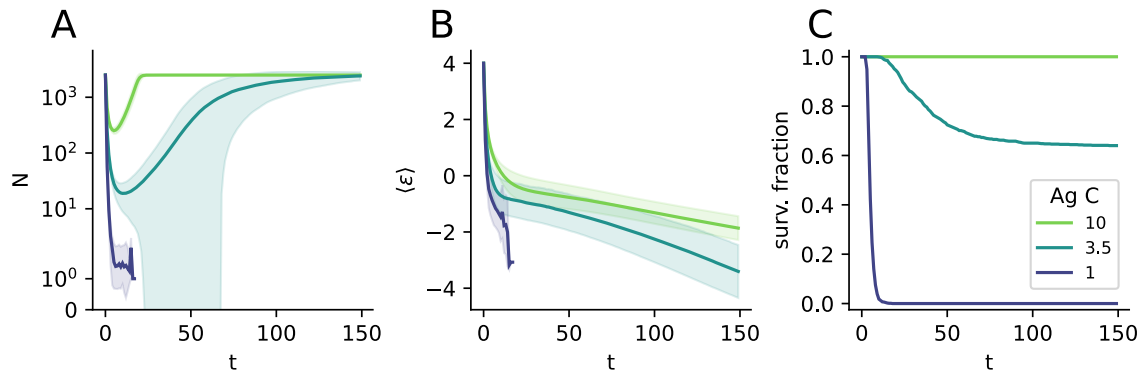


FIG. 3. Average evolution of 1000 different stochastic simulations of the model at three different levels of Ag concentration  $C = 1, 3.5, 10$ , color-coded according to the legend on the right. **A**: population size  $N$  as a function of evolution round. Shaded area covers one standard deviation for surviving simulations. The minimum population size on the bottleneck depends strongly on Ag concentration **B**: same as panel A but for the average population binding energy  $\langle \epsilon \rangle$ . Notice how for surviving populations the maturation speed depends on Ag concentration. **C**: Fraction of surviving simulations as a function of time. At low concentration the bottleneck drives all simulations to extinction, while at high concentration the population survives with high probability.

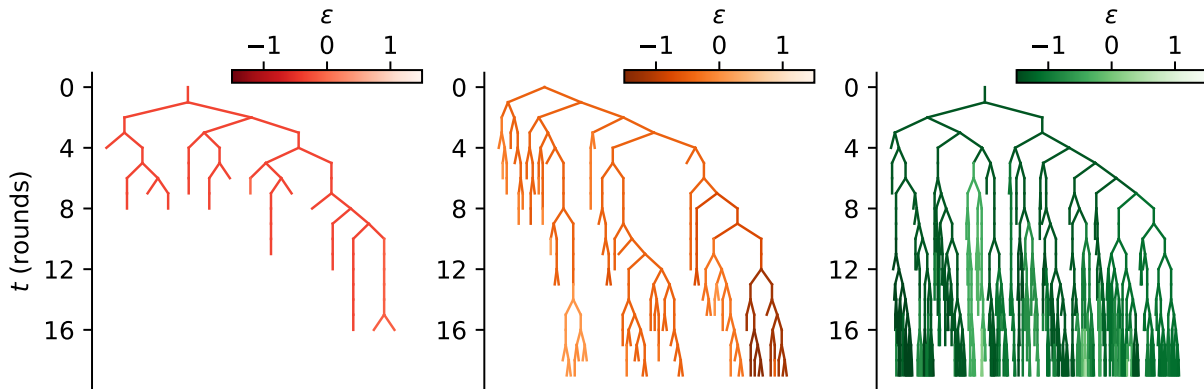


FIG. 4. Examples of stochastic lineage evolution through a population bottleneck. We perform a single simulation of our model at Ag concentration  $C = 5$  and consider three different progenitors with different initial affinities (red  $\epsilon_i \sim -0.3$ , orange  $\epsilon_i \sim -0.45$  and green  $\epsilon_i \sim -1.3$ ). We represent their progeny evolution in the form of a tree with each cell corresponding to a node, and encoding affinity in the branch color. The lineage of the red progenitor quickly goes extinct, while the lineage of the orange progenitor survives the bottleneck but only with few individuals. The green progenitor lineage conversely survives the population bottleneck and undergoes great expansion. Notice how fate correlates with the initial progenitor affinity.

daughter cell this probability is more easily expressed as one minus the probability of not being removed:

$$d_1(\epsilon) = \left[ 1 - \int d\Delta \epsilon K(\Delta \epsilon) P_S(\epsilon + \Delta \epsilon) (1 - p_{\text{diff}}) \right]^2, \quad (4)$$

where the expression for  $K(\Delta)$  is the one given in eq. (1), and  $P_S(\epsilon)$  is the probability for a cell with binding energy  $\epsilon$  of surviving selection. In the bottleneck state most of the selection pressure is generated by Ag-binding selection (i.e.  $\bar{\epsilon}_t < \epsilon_{\text{Ag}}$ ). As a first approximation we therefore neglect competitive selection for T-cell help, and consider simply  $P_S(\epsilon) = P_{\text{Ag}}(\epsilon)$  (cf. eq. (2)). This introduces two important simplifications. First, the expression of  $P_S$  does not depend on time. Second, removing the competitive selection decouples the fate of all cells in the

population.

The probabilities  $d_t(\epsilon)$  for  $t > 1$  can be evaluated using recursive relations that express the probability of extinction in  $t$  rounds as the probability for each daughter cell to either go extinct in one round, or to survive the first round but to have their respective offspring go extinct in the remaining  $t - 1$  rounds:

$$d_t(\epsilon) = \left[ 1 - \int d\Delta \epsilon K(\Delta \epsilon) P_S(\epsilon + \Delta \epsilon) (1 - p_{\text{diff}}) \times (1 - d_{t-1}(\epsilon + \Delta \epsilon)) \right]^2 \quad (5)$$

In other words, the probability that all of the offspring goes extinct in  $t$  rounds is the probability that each of

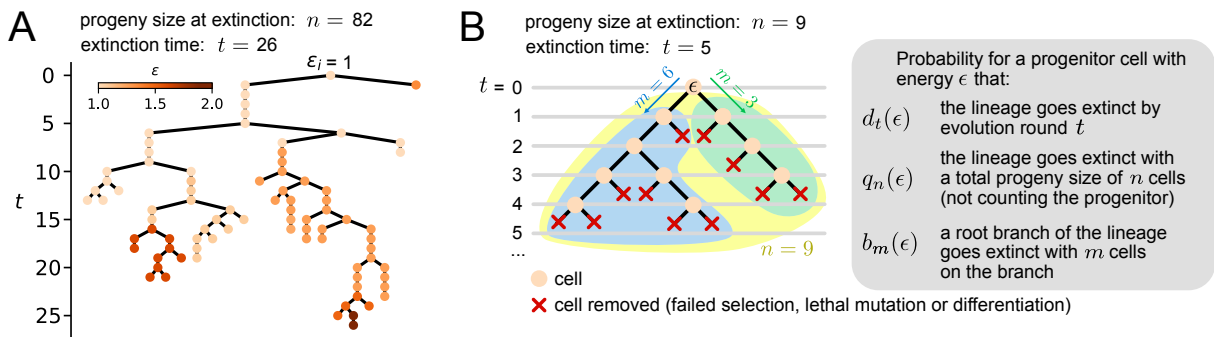


FIG. 5. **A.** Example of lineage issued from a progenitor with binding energy  $\epsilon_i = 1$  obtained from a stochastic simulation performed at Ag concentration  $C = 7$  in the approximation of only Ag-binding selection. Each node in the tree represents a cell, its binding energy  $\epsilon$  encoded using the colorscale on top. In this example cells progressively accumulate deleterious mutations until after 26 evolution rounds the lineage eventually goes extinct. **B.** Schematic illustration of the quantities analyzed in our theory. On the left we depict a lineage evolution, stemming from a progenitor with binding energy  $\epsilon$ . The probability that such a lineage goes extinct by time  $t$  is indicated with  $d_t(\epsilon)$  (in this example  $t = 5$ ). The quantity  $q_n(\epsilon)$  represents instead the probability that the lineage goes extinct counting a total of  $n$  cells not including the progenitor ( $n = 9$  in this example). Finally,  $b_m(\epsilon)$  is the probability that one of the two sub-lineage stemming from the progenitor goes extinct counting  $m$  cells (here  $m = 6$  for the left branch and  $m = 3$  for the right one).

the lineages stemming from the two daughter cells generated during the duplication phase of the first round are removed before the end of round  $t$ . Since division is symmetric this probability must be the same for each daughter cell, and is the term inside the square brackets. This term is more easily expressed as one minus the probability that the lineage survives (term in the integral). In turn this can be decomposed as the probability that the daughter cell survives mutation (possibly with an energy change of entity  $\Delta\epsilon$ ), selection and differentiation, multiplied by the probability that its offspring does not go extinct by the following  $t-1$  rounds (term on the second line). A more extended explanation for the derivation of this equation and its numerical evaluation is provided in appendix C.

In fig. 6A we plot the behavior of  $d_t(\epsilon)$  as a function of evolution round  $t$  and progenitor binding energy  $\epsilon$  (orange curves, color indicates extinction round  $t$ ). As expected, the extinction probability is a monotonically increasing function of time and of energy, and reaches an asymptotic value  $d_\infty(\epsilon)$  for large  $t$ . Our analytical result is in excellent agreement with simulations for the mean extinction probability (blue dots). The asymptotic probability  $d_\infty(\epsilon)$  ranges between:

- $d_\infty(\epsilon \rightarrow +\infty) = 1$ . This is easily understood, since high-energy (i.e. low-affinity) cells will not pass the selection step and their progeny will quickly go extinct.
- $d_\infty(\epsilon \rightarrow -\infty) = \min\{1, [1 - 1/\alpha]^2\}$ , with  $\alpha = (1 - p_{\text{let}})(1 - p_{\text{diff}})$ . The value of extinction probability  $d_\infty$  for very high-affinity cells may be higher than 0 since lethal mutations and differentiation may still drive the lineage to extinction, especially during the first few evolution rounds when the offspring size is still small. The above expression

for  $d_\infty(\epsilon \rightarrow -\infty)$  can be obtained by searching a fixed point to eq. (5), and considering that, for  $\epsilon \rightarrow -\infty$ , mutations do not sensibly change the survival probability and can be neglected. The parameter  $\alpha$  defined above is then the probability for a high-affinity cell to survive one round and not be removed by lethal mutations or through differentiation. Notice that, if  $\alpha < 1/2$ , we have  $d_\infty(\epsilon \rightarrow +\infty) = 1$ , as is to be expected when on average less than one individual in the offspring will survive. However this case is pathological: in this regime, irrespective of the progenitor energy, the population always goes extinct ( $d_\infty(\epsilon) = 1$ ).

It is worthwhile to notice that the “infinite-time” extinction probability  $d_\infty$  described in our theory represents the probability for a lineage to go extinct in the bottleneck phase, and does not reflect the probability of lineage fixation or extinction when the population is at maximum size. In eqs. (4) and (5) we indeed neglected the carrying capacity constraint. This is justified when considering the bottleneck state, in which the population is not at its maximum size. However if some lineages are able to survive this state, then the population will eventually grow back to maximum size, and some cells will be randomly removed through enforcement of the carrying capacity. One should therefore extend the theory and add another term to correctly describe the lineage survival probability far away from the bottleneck.

## 2. Distribution of extinction times

The probability that a lineage generated by a progenitor with energy  $\epsilon$  goes extinct exactly at round  $t$  can

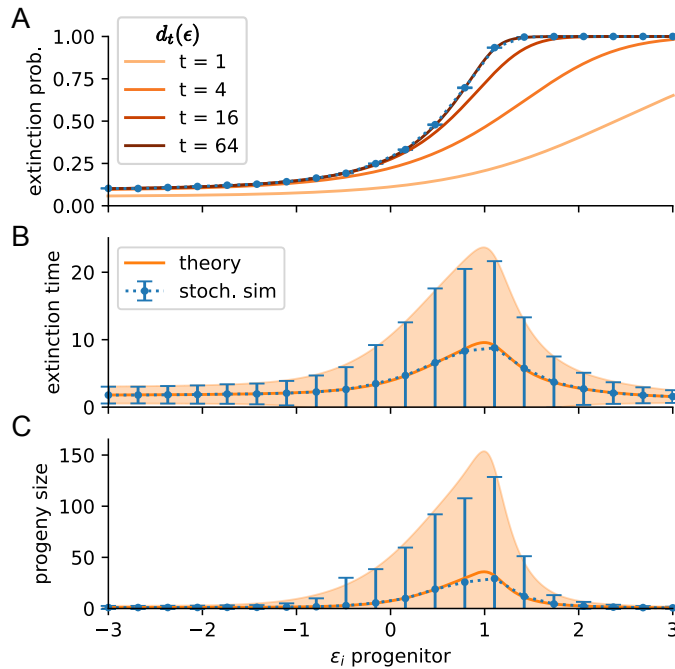


FIG. 6. Comparison between stochastic simulations (blue) and theory (orange) for the probability of extinction (A), lineage extinction time (B) and average progeny size at extinction (C) as a function of the progenitor energy  $\epsilon_i$  in absence of competitive selection. For each conditions we consider 5000 different stochastic simulations that terminate with extinction at Ag concentration  $C = 7$ . **A:** stochastic extinction probability (blue dots, error bar indicate the standard error of the mean) evaluated as the fraction of simulations that terminate with extinction over the total number of simulations performed. This is compared to the value of  $d_t(\epsilon)$  as described by our theory. **B:** mean and standard deviation of extinction time (blue) over 5000 simulations terminating in extinction. This is compared to the theoretical prediction (orange) for the mean and standard deviation of this quantity, obtained using the time extinction probability  $r_t(\epsilon)$ . **C:** same as B but for the progeny size. In this case the theoretical predictions are obtained using the generating function theory.

easily be expressed as

$$r_t(\epsilon) = d_t(\epsilon) - d_{t-1}(\epsilon). \quad (6)$$

This allows us to evaluate the mean and variance for the extinction time probabilities (see fig. 6B) simply from the first two moments of the distribution:

$$\langle t \rangle_\epsilon = \sum_{t=0}^{\infty} t r_t(\epsilon), \quad \langle t^2 \rangle_\epsilon = \sum_{t=0}^{\infty} t^2 r_t(\epsilon), \quad (7)$$

In fig. 6B we compare, in the approximation of no competitive selection, the average extinction time computed from simulations (blue, error bars indicate the standard deviation of extinction times for each progenitor affinity) with theoretical predictions (orange, shaded area covers one standard deviations). We again find a very good match. The average extinction time shows a peak for intermediate affinities, which can be interpreted as follows. Low-affinity progenitors, i.e. having high binding energy have close-to-one probability of extinction, and very often go extinct in the first few rounds. High-affinity cells on the contrary have a small but non-zero probability of extinction, see value of extinction probability in fig. 6A. This is mainly due to affinity-independent terms (differentiation and lethal mutation probabilities) which confer

to the lineage a small chance of going extinct during the first few evolution rounds, when the progeny size is still small. For affinities close to  $\epsilon = 1$  we observe intermediate values for the probability of survival, and maximum value for average extinction time. This behaviour can be better understood when mutations are turned off, in which case equations can be solved exactly, as shown below.

### 3. Exactly solvable case of no mutation

We hereafter consider the case of no affinity-affecting mutations, for which the mutation kernel eq. (1) reads

$$K(\Delta) = (1 - p_{\text{let}}) \delta(\Delta). \quad (8)$$

In this case genealogies belong to the class of Galton-Watson trees [24], and the asymptotic survival probability can be derived exactly. This probability is better expressed by considering the quantity

$$\gamma(\epsilon) = (1 - p_{\text{let}}) P_S(\epsilon) (1 - p_{\text{diff}}) \quad (9)$$

which represents the probability for a daughter cell to remain in the GC and not be removed by either lethal



mutations, selection or differentiation. The infinite-time extinction probability  $d_\infty(\epsilon)$  can be found by rewriting eq. (5) in the limit  $t \rightarrow \infty$ :

$$d_\infty(\epsilon) = \min \left( 1, \left( \frac{1}{\gamma(\epsilon)} - 1 \right)^2 \right). \quad (10)$$

As expected, lineages will always go extinct if the average number of surviving offspring at division is not greater than one:  $d_\infty(\epsilon) = 1$  if  $\gamma(\epsilon) \leq \frac{1}{2}$ . In fig. 7A we report the behavior of  $d_\infty$  as a function of  $\gamma$ . This function presents a singularity at the critical value  $\gamma = 1/2$ , for which the Galton-Watson process is critical.

Finding an explicit expression for the distribution of extinction times is harder, but results can be obtained for the critical value  $\gamma = \frac{1}{2}$ . We find that the extinction time probability behaves asymptotically as a power law, with infinite mean and variance:  $r_t \sim 4/t^2$  for large  $t$ . This result, which is a known feature of critical Galton-Watson processes, can be verified by inserting the Ansatz  $d_t \sim 1 - \alpha t^{-1} + o(t^2)$  in eq. (5), together with the simplified form of the mutation kernel (8) and the assumption that  $\gamma = 1/2$ . The only admissible solution is  $\alpha = 4$  which, combined with the definition of  $r_t$  eq. (6), proves the statement. In fig. 7B we display the mean and variance of the extinction time distribution as a function of  $\gamma$ . Comparison with fig. 6B shows that the divergence is removed when evolution includes affinity-affecting mutations. Mutations drive lineages away from the critical line, either to high affinities and survival, or to lower affinities and extinction.

### B. Case of full population

Building on the results derived above, we now turn to the problem of quantifying the average probability of extinction for the whole population.

As a first approximation we do not consider competitive selection, since most of the selection pressure in a bottleneck is given by Ag-binding selection. Given a total of  $N_i$  cells in the initial population, having energies  $\{\epsilon_k\}_{k=1\dots N_i}$  the probability that all cells will be extinct by evolution round  $t$  is simply given by the product of extinction probabilities for all cells  $\prod_k d_t(\epsilon_k)$ . Moreover, since the initial energies are independently extracted from a Gaussian distribution  $\varphi(\epsilon)$  with mean  $\mu_i$  and standard deviation  $\sigma_i$ , the average extinction probability by round  $t$  over all possible extractions of the initial population is given by:

$$P_{\text{ext}}(t) = \left( \int d\epsilon \varphi(\epsilon) d_t(\epsilon) \right)^{N_i} \quad (11)$$

With the help of this formula we evaluate the average survival probability as a function of Ag concentration  $C$  and initial population size  $N_i$ , and compare the prediction with stochastic simulations in which we turn off

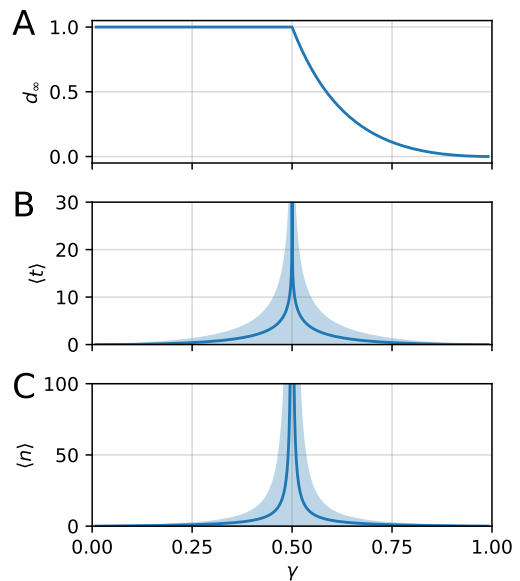


FIG. 7. Value of the extinction probability  $d_\infty$  (A), average extinction time  $\langle t \rangle$  (B) and average progeny size  $\langle n \rangle$  (C) as a function of the survival probability  $\gamma(\epsilon)$  (cf. eq. (9)) in the approximation of no affinity-affecting mutation. In B and C shaded area covers one standard deviation. Notice how extinction times and genealogy sizes diverge at  $\gamma = \frac{1}{2}$ .

T-cell selection. The results, reported in fig. 8B and C (blue), match exactly.

In the presence of competitive selection the empirical survival probability evaluated from simulations slightly decreases, compare blue and orange dotted lines in fig. 8B and C. The theory can be extended to account for T-selection in an effective manner. In practice, one needs first to extend the theory to include a time-dependence of the survival probability. At this point competitive selection can be included by introducing an effective coupling between cells in a ‘mean field’ fashion, by estimating the average evolution of the term  $\bar{\epsilon} = -\log\langle e^{-\epsilon} \rangle_{\text{pop}}$  contained in the expression for the T-selection survival probability eq. (3).

Assume that the survival probability  $P_S(\epsilon, t)$  is now time-dependent. The probability of extinction does not depend solely on the number of evolution rounds anymore, but also on the initial time at which the progenitor is considered. We define  $d_{t,s}(\epsilon)$  as the probability that a cell, which at the end of round  $t$  has binding energy  $\epsilon$ , will have all of its offspring extinct by the end of round  $s > t$ . For any value  $t \geq 0$  we can write as before the probability of extinction in one round:

$$d_{t,t+1}(\epsilon) = \left( 1 - \int d\Delta K(\Delta) P_S(\epsilon + \Delta, t) (1 - p_{\text{diff}}) \right)^2 \quad (12)$$

And for any pair of rounds  $s > t \geq 0$ , with  $s - t > 1$ , the

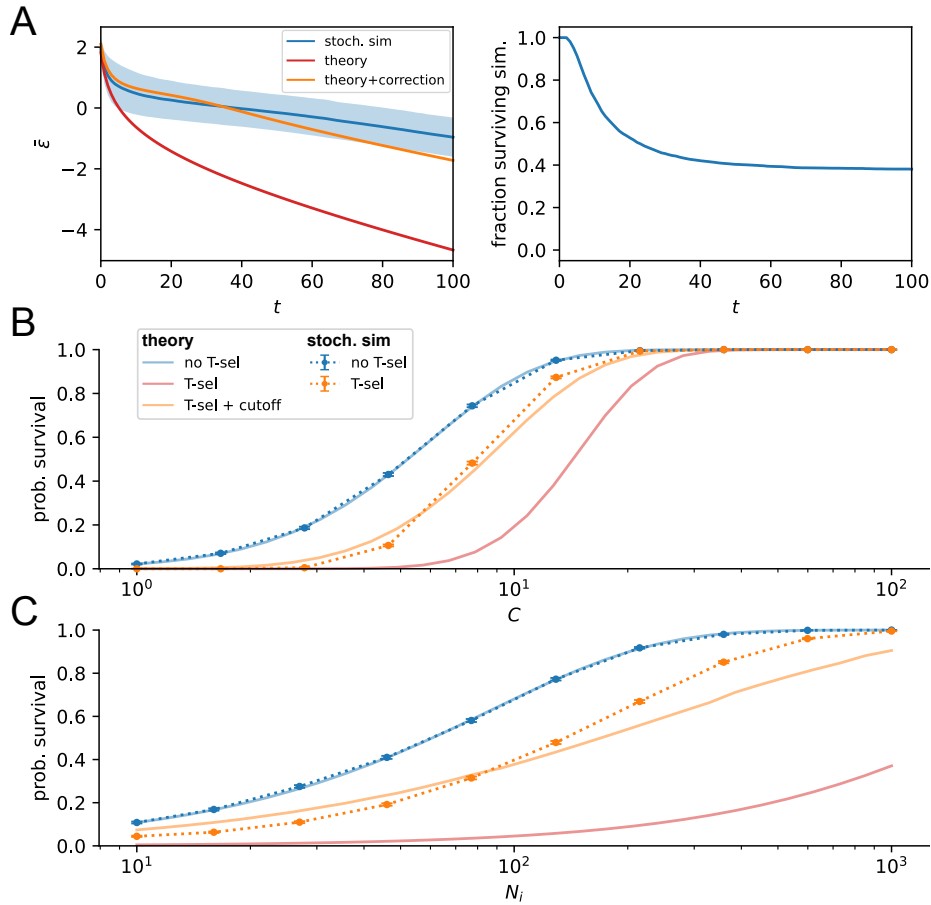


FIG. 8. Probability of population survival in a bottleneck condition as a function of initial population size  $N_i$  and Ag concentration  $C$ . **A:** Left: comparison between the evolution of  $\bar{\epsilon}$  in stochastic simulations (blue, mean and standard deviation over 5000 simulations) and theoretical prediction without (red) and with (orange) finite-size correction. This correction consists in cutting the tail of the initial energy distribution in proximity of the expected value for the highest-affinity individual. The correction improves the prediction for the evolution of  $\bar{\epsilon}$  at short times. Right: fraction of surviving simulations as a function of evolution round. With the finite-size correction the value of  $\bar{\epsilon}$  is well-approximated during the time it takes for most of the simulations to go extinct. In this example we set  $C = 7$ ,  $N_i = 100$ . **B:** bottleneck survival probability as a function of antigen concentration  $C$ . Comparison between stochastic simulations (dotted line, error bars indicate the standard error of the mean) and the predictions our theory (full lines). Stochastic simulations are reported both without (blue) and with (orange) competitive selection for T-cell help (T-sel). For the theory instead we consider the case without T-sel (blue), with T-sel (red) and with T-sel plus finite-size correction (orange). In the absence of T-sel all cells evolve independently, and the theory and simulations match exactly. The inclusion of T-sel slightly decreases the survival probability in stochastic simulations. Accounting for this contribution by using the infinite-size estimate for the evolution of  $\bar{\epsilon}$  overestimates the selection pressure. Adding the finite-size correction results in a much better estimate. In this example we set  $N_i = 100$ . **C:** same as B, but the survival probability is evaluated as a function of the initial size of the population  $N_i$ . Here we set  $C = 7$ .

following recursive relation, analogous to eq. (5), holds:

$$d_{t,s}(\epsilon) = \left[ 1 - \int d\Delta K(\Delta) P_S(\epsilon + \Delta, t) (1 - p_{\text{diff}}) \times (1 - d_{t+1,s}(\epsilon + \Delta)) \right]^2 \quad (13)$$

Finally, similarly to eq. (11), the probability that the full population goes extinct by evolution round  $t$  is given by:

$$P_{\text{ext}}(t) = \left( \int d\epsilon \varphi(\epsilon) d_{0,t}(\epsilon) \right)^{N_i}, \quad (14)$$

where  $\varphi(\epsilon)$  is a Gaussian distribution with mean  $\mu_i$  and standard deviation  $\sigma_i$ .

At this point we can make explicit the time dependence of the survival probability including selection for T-cell help:  $P_S(\epsilon, t) = P_{\text{Ag}}(\epsilon) P_{\text{T}}(\epsilon, \bar{\epsilon}_t)$  (cf. eq. (3)). Given the stochastic nature of our model, the variable  $\bar{\epsilon}_t$  which quantifies selection pressure is in reality a stochastic variable. We estimate its average evolution using the large-population-size limit described in appendix B, under which the model becomes deterministic. This allows us to numerically evaluate the extinction probabil-

ity eq. (14). The outcome, however, underestimates the real survival probability (compare red curve and orange dotted line in fig. 8B and C). This mismatch originates mainly from the fact that in the big-size approximation  $\bar{\epsilon}$  evolves faster than in stochastic simulations (cf. blue and orange line in fig. 8A-left). In turn, this occurs because the value of  $\bar{\epsilon}$  is strongly dependent on the high-affinity tail of the population, whose evolution is influenced by finite-size effects.

This discrepancy can, however, be reduced with a simple finite-size correction. This correction is based on the consideration that the large-size limit of the model approximates the population binding energy histogram with a continuous distribution, encoded in the density function  $\rho_t(\epsilon)$  (cf. appendix B). At the beginning of evolution this function takes the shape of a normal distribution, corresponding to the initial binding energy distribution of naive responders, with tails extending indefinitely in both directions. As the population is finite in reality, consisting of  $N_i$  individuals, we do not expect these tails to be populated. The correction procedure consists in removing these tails, by setting the initial distribution equal to zero outside a range delimited by two values  $[\epsilon^-, \epsilon^+]$ .

These two values are chosen equal to the expected energy of, respectively, the highest and lowest affinity individual in the population. The probability distribution for their binding energies can be expressed as a function of the naive binding energy distribution  $\varphi(\epsilon)$  (as before a Gaussian with mean  $\mu_i$  and variance  $\sigma_i^2$ ) from which the energy of all cells is extracted. If we call  $F(\epsilon) = \int_{-\infty}^{\epsilon} d\epsilon' \varphi(\epsilon')$  the cumulative distribution function, then these distributions can be expressed as:

$$\varphi^+(\epsilon) = \frac{d}{d\epsilon} [F(\epsilon)]^{N_i} \quad (15)$$

$$\varphi^-(\epsilon) = -\frac{d}{d\epsilon} [1 - F(\epsilon)]^{N_i} \quad (16)$$

The values  $\epsilon^\pm$  simply correspond to the means of these distributions.

Removing the tails to the initial distribution causes an initial slow-down in the evolution of  $\bar{\epsilon}$  (cf. green line in fig. 8A-left). This slow-down is eventually lost, but the agreement remains for a time sufficient for most of the stochastic simulations to go extinct (cf. fig. 8A-right) which is the relevant timescale to capture bottleneck survival.

Taking the value of  $\bar{\epsilon}$  obtained by combining the big-size approximation (cf. appendix B) with the cutoff correction described above, and using it to evaluate the population survival probability, we obtain a much better agreement of the theory with simulations (compare orange curve and orange dotted line in fig. 8B and C). The remaining discrepancy are due to the fact that the average evolution of  $\bar{\epsilon}$  is still not exactly captured, and the ‘mean-field’ nature of our approximation, which neglects the feedback of the energies in the population onto  $\bar{\epsilon}$ .

#### IV. LINEAGE SIZE AT EXTINCTION

In this Section we focus on the distribution of sizes of the progeny at extinction. This size strongly depends on the model parameters, such as the energy of the progenitor. An example is displayed in fig. 5A, in which the lineage consists of a total of 82 cells. Like extinction time, this quantity is well-defined only for lineages that go extinct. Populations that are able to pass the bottleneck undergo exponential growth, with a rate that can be calculated from the large-size theory of Appendix B, see [13].

##### 1. Recursion equations for the distribution of sizes

Similarly to what was done in the previous Section for the extinction time and probability, we now derive a recursive formula to quantify the total offspring size. We need to keep track of the sum of two random variables representing the numbers of descendants of each daughter cell. The recursion therefore includes a convolution, which is numerically harder to compute but can be handled using probability generating functions. The recursive relations can be expressed in term of these functions, and can be used to evaluate the moments of the probability distribution without having to numerically perform the convolution.

We name  $q_n(\epsilon)$  the probability that a progenitor with energy  $\epsilon$  generates a total offspring of exactly  $n$  cells before extinction, not counting the progenitor itself (see fig. 5B). This probability can be better expressed if we separate the contribution of the two daughter cells. Considering genealogies encoded as binary trees, we call  $b_m(\epsilon)$  the probability that along the branch corresponding to one of the daughter cells of a progenitor with energy  $\epsilon$  we find a total of  $m$  descendants (including the daughter cell itself) before extinction (see fig. 5). The expression for  $m = 0$  is simply given by the probability that the daughter cell is removed before the end of the round:

$$\begin{aligned} b_0(\epsilon) &= 1 - \int d\Delta K(\Delta) P_S(\epsilon + \Delta) (1 - p_{\text{diff}}) \\ &= \sqrt{d_1(\epsilon)} \end{aligned} \quad (17)$$

The recursive relation in this case is composed of two equations. The first is a convolution that decomposes the probability of having  $n$  descendants as a sum over all possible repartitions of the descendant number along the two branches:

$$q_n(\epsilon) = \sum_{m=0}^n b_m(\epsilon) b_{n-m}(\epsilon) \quad (18)$$

The second expresses the probability to find  $m$  descendants along a branch as the probability that the daughter

cell survives and has  $m - 1$  descendants:

$$b_m(\epsilon) = \int d\Delta K(\Delta) P_S(\epsilon + \Delta) (1 - p_{\text{diff}}) q_{m-1}(\epsilon + \Delta) \quad (19)$$

We introduce the generating functions  $Q(z, \epsilon)$  and  $B(z, \epsilon)$ , defined as:

$$Q(z, \epsilon) = \sum_{n=0}^{\infty} q_n(\epsilon) z^n, \quad B(z, \epsilon) = \sum_{m=0}^{\infty} b_m(\epsilon) z^m. \quad (20)$$

In terms of these generating functions equations eqs. (18) and (19) become

$$Q(z, \epsilon) = B(z, \epsilon)^2 \quad (21)$$

$$\frac{1}{z} [B(z, \epsilon) - b_0(\epsilon)] = \int d\Delta K(\Delta) P_S(\epsilon + \Delta) \times (1 - p_{\text{diff}}) Q(z, \epsilon + \Delta) \quad (22)$$

These relations can be used to evaluate the moments of these distributions with two additional considerations. The first is that  $\sum_{n=0}^{\infty} q_n(\epsilon) = d_{\infty}(\epsilon)$ . This sum does not converge to one since it only considers lineages that eventually go extinct. For the functions  $Q$  and  $B$  this translates into:

$$Q(z = 1, \epsilon) = d_{\infty}(\epsilon), \quad B(z = 1, \epsilon) = \sqrt{d_{\infty}(\epsilon)}. \quad (23)$$

Secondly, the moments of the distributions can be evaluated from the generating functions as:

$$\begin{aligned} \langle n^k \rangle_{\epsilon} &= \frac{1}{d_{\infty}(\epsilon)} \sum_{n=0}^{\infty} n^k q_n(\epsilon) \\ &= \frac{1}{d_{\infty}(\epsilon)} (z \partial_z)^k Q(z, \epsilon) |_{z=1} \end{aligned} \quad (24)$$

$$\begin{aligned} \langle m^k \rangle_{\epsilon} &= \frac{1}{\sqrt{d_{\infty}(\epsilon)}} \sum_{m=0}^{\infty} m^k b_m(\epsilon) \\ &= \frac{1}{\sqrt{d_{\infty}(\epsilon)}} (z \partial_z)^k B(z, \epsilon) |_{z=1} \end{aligned} \quad (25)$$

Applying the operator  $z \partial_z$  one and two times on eq. (21) restitutes the following relations between the first two moments:

$$\langle n \rangle_{\epsilon} = 2 \langle m \rangle_{\epsilon}, \quad \langle n^2 \rangle_{\epsilon} = 2 \langle m^2 \rangle_{\epsilon} + 2 \langle m \rangle_{\epsilon}^2 \quad (26)$$

This corresponds simply to the fact that the total number of descendants is the sum of the descendants along the two branches. Applying the same operator on eq. (22) gives:

$$\sqrt{d_{\infty}(\epsilon)} \langle m \rangle_{\epsilon} = \int d\Delta K(\Delta) P_S(\epsilon + \Delta) (1 - p_{\text{diff}}) \times d_{\infty}(\epsilon + \Delta) [2 \langle m \rangle_{\epsilon + \Delta} + 1] \quad (27)$$

$$\begin{aligned} \sqrt{d_{\infty}(\epsilon)} \langle (m-1)^2 \rangle_{\epsilon} &= \int d\Delta K(\Delta) P_S(\epsilon + \Delta) \\ &\times (1 - p_{\text{diff}}) d_{\infty}(\epsilon + \Delta) [2 \langle m^2 \rangle_{\epsilon + \Delta} + 2 \langle m \rangle_{\epsilon + \Delta}^2 + 1] \end{aligned} \quad (28)$$

These equations can be solved numerically if we express them as fixed-point equations for the functions  $\langle m \rangle_{\epsilon}$  and  $\langle m^2 \rangle_{\epsilon}$ :

$$\langle m \rangle_{\epsilon} = \frac{1}{\sqrt{d_{\infty}(\epsilon)}} \int d\Delta K(\Delta) P_S(\epsilon + \Delta) (1 - p_{\text{diff}}) \times d_{\infty}(\epsilon + \Delta) [2 \langle m \rangle_{\epsilon + \Delta} + 1] \quad (29)$$

$$\begin{aligned} \langle m^2 \rangle_{\epsilon} &= \frac{1}{\sqrt{d_{\infty}(\epsilon)}} \int d\Delta K(\Delta) P_S(\epsilon + \Delta) (1 - p_{\text{diff}}) \\ &\times d_{\infty}(\epsilon + \Delta) [2 \langle m^2 \rangle_{\epsilon + \Delta} + 2 \langle m \rangle_{\epsilon + \Delta}^2 + 4 \langle m \rangle_{\epsilon + \Delta} + 1] \end{aligned} \quad (30)$$

The moments for  $n$  can then easily be evaluated using eq. (26).

In fig. 6C we compare the theoretical prediction for the first two moments (orange line represents the mean and shaded area covers one standard deviation) with the corresponding quantities from stochastic simulations (blue, error bars cover one standard deviation). Once more we find a good match. The peak at intermediate affinities can be explained, as done above for the extinction time, considering the critical nature of this phenomenon at intermediate values of the binding energy. This is done in the next section.

## 2. Exactly solvable case of no mutation

Similarly to what done for extinction probability, in the absence of affinity-affecting mutations we can find an explicit expression for the mean and variance of the population size at extinction. By plugging the simplified expression for the mutation kernel eq. (8) into eqs. (21) and (22) one obtains the following second degree equation for the generating function  $B$ :

$$z \gamma(\epsilon) B(z, \epsilon)^2 - B(z, \epsilon) + 1 - \gamma(\epsilon) = 0 \quad (31)$$

Where as before  $\gamma(\epsilon)$  is the probability for a daughter cell not to be removed from the population during the evolution round, cf. eq. (9). This equation has two solutions. The correct one can be chosen by considering that  $B$  must be a monotonically increasing function of  $z$ . This gives:

$$B(z, \epsilon) = \frac{1 - \sqrt{1 - 4z\gamma(\epsilon)(1 - \gamma(\epsilon))}}{2z\gamma(\epsilon)} \quad (32)$$

The function  $Q(z, \epsilon)$  can be evaluated from eq. (21), and the mean and variance for the population extinction sizes can be obtained using eq. (24). This results in:

$$\langle n \rangle_{\epsilon} = \begin{cases} \frac{2\gamma(\epsilon)}{1-2\gamma(\epsilon)} & \text{if } \gamma(\epsilon) < 1/2 \\ \frac{2-2\gamma(\epsilon)}{2\gamma(\epsilon)-1} & \text{if } \gamma(\epsilon) > 1/2 \end{cases} \quad (33)$$

$$\langle n^2 \rangle_{\epsilon} - \langle n \rangle_{\epsilon}^2 = \langle n \rangle_{\epsilon} (\langle n \rangle_{\epsilon} + 1) (\langle n \rangle_{\epsilon} + 2) / 2 \quad (34)$$

This quantity is reported as a function of  $\gamma$  in fig. 7C. Similarly to what observed for the mean and variance of

the extinction time, Both of these quantities diverge for the critical value  $\gamma = \frac{1}{2}$ , but this divergence is removed when mutations are considered.

It is interesting to consider the effect of this divergence on the coefficients  $q_n(\epsilon)$ , that represent the probability of a lineage that stems from a progenitor with binding energy  $\epsilon$  to go extinct with a total progeny of  $n$  cells (see fig. 5B). From the definition of the generating function  $Q(z, \epsilon)$  (cf. eq. (20)) it follows that the coefficients  $q_n$  can be obtained by Taylor expansion of this function around  $z = 0$ . In turn  $Q$  can be easily obtained from eq. (32) using the property  $Q = B^2$  (see eq. (21)). The expansion results in the following expression for the coefficients when mutations are absent:

$$q_n = \frac{(2n+2)!}{(n+1)!(n+2)!} \gamma^n (1-\gamma)^{n+2} \quad (35)$$

$$\stackrel{n \gg 1}{\approx} \frac{(1-\gamma)^2}{\sqrt{\pi n^3}} (4\gamma(1-\gamma))^n$$

In general these probabilities decay exponentially fast as a function of the size  $n$ . At the critical value  $\gamma = 1/2$  however the term  $4\gamma(1-\gamma)$  becomes equal to 1 and the coefficients algebraically decay as  $n^{-3/2}$ . The asymptotic decay of the coefficients  $q_n$  for  $n \gg 1$  can also be obtained from the behavior of the generating function  $Q$  around its singularity at  $z_c = 1/(4\gamma(1-\gamma))$ . In particular  $Q \propto (1 - z/z_c)^{1/2}$  close to the singularity, and therefore  $q_n \propto n^{-3/2} z_c^{-m}$  [25], which gives the expected asymptotic behavior described above.

### 3. Case of small-effect mutations

Corrections to the asymptotic behavior above arise when small-effect mutations are considered. In particular in eq. (8) we substitute the Dirac delta distribution with a peaked Gaussian, and consider a mutation kernel having the form:

$$K(\Delta) = (1 - p_{\text{let}}) \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\Delta^2}{2\sigma^2}\right) \quad \text{with } \sigma \ll 1 \quad (36)$$

When the standard deviation is small enough one can approximate the integrals in eqs. (17) and (22) by Taylor-expanding the functions that multiply the mutation kernel around  $\Delta = 0$ . This results in the following approximation for eq. (22):

$$B - 1 + \gamma + \frac{\sigma^2}{2} \gamma'' = z\gamma B^2 + z \frac{\sigma^2}{2} [\gamma B^2]'' \quad (37)$$

according to the definition of  $\gamma$  (cf. eq. (9)), and with inverted commas indicating derivatives with respect to  $\epsilon$ . This equation is analogous to eq. (31) with the addition of perturbation terms of the order of  $\sigma^2$ . We therefore suppose  $B(z, \epsilon) \sim B_0(z, \epsilon) + \sigma^2 \Delta B(z, \epsilon)$  where  $B_0$  is given by eq. (32) and  $\Delta B$  represents a perturbation to this  $\sigma^2 = 0$  solution. By plugging this Ansatz in the

previous equation we find the following expression for the perturbation:

$$\Delta B = \frac{(B_0 - 1)}{2\gamma r^4} \left[ 2z\gamma'^2(1-2\gamma) + r \frac{\gamma'^2}{\gamma} + r^2 \left( \gamma'' - \frac{\gamma'^2}{\gamma} \right) \right], \quad (38)$$

where  $r$  is defined through

$$r = \sqrt{1 - 4z\gamma(1-\gamma)} = \sqrt{\frac{z_c - z}{z_c}}. \quad (39)$$

The critical value of  $z$  is given as before by  $z_c = 1/(4\gamma(1-\gamma))$ . Moreover, when the survival probability  $P_S$  is given by eq. (2), the derivatives of  $\gamma$  (cf. eq. (9)) can be expressed as

$$\gamma' = -\gamma(1-\tilde{\gamma}) \quad (40)$$

$$\gamma'' = \gamma(1-\tilde{\gamma})(1-2\tilde{\gamma}), \quad (41)$$

where  $\tilde{\gamma} = P_{\text{Ag}}(\epsilon)$ . From the expression for  $\Delta B$  we can derive the perturbation to  $Q$  by using eq. (21). Keeping only the higher order in  $\sigma^2$ , we obtain

$$Q \sim Q_0 + \sigma^2 \Delta Q, \quad \text{with } \Delta Q = 2 B_0 \Delta B \quad (42)$$

As stated in the previous section, the behavior of the probabilities  $q_n$  is strictly related to the behavior of  $Q$  around its singularity  $z_c$ . In particular the magnitude of the perturbation to the coefficients  $q_n$  introduced by mutations can be derived from the study of  $\Delta Q$ . If we operate the substitution  $z = z_c(1-r^2)$  and expand  $\Delta Q$  in powers of  $r$  we obtain:

$$\Delta Q = c_4 r^{-4} + c_3 r^{-3} + c_2 r^{-2} + c_1 r^{-1} + O(1) \quad (43)$$

with the following expressions for the coefficients:

$$\begin{aligned} c_4 &= (1-2\gamma)^2(1-\tilde{\gamma})^2 \\ c_3 &= -(1-2\gamma)^2(1-\tilde{\gamma})^2 \\ c_2 &= -2(1-\gamma)(1-\tilde{\gamma})(1-2\gamma\tilde{\gamma}) \\ c_1 &= 2(1-\gamma)(1-\tilde{\gamma})(3-2\gamma-2\gamma\tilde{\gamma}) \end{aligned} \quad (44)$$

Based on this expansion we can derive an expression for the perturbation to the coefficients  $\Delta q_n = q_n - q_n^0$  caused by weak mutations [25], where  $q_n^0$  represents the value in the absence of mutations (cf. section IV 2). We obtain that, for large  $n$ ,

$$\begin{aligned} \Delta q_n &= \sigma^2 [c_4(n+1) + c_3(n/\pi)^{1/2}(2+3/4n) + c_2 \\ &\quad + c_1(\pi n)^{-1/2} + O(n^{-3/2})] (4\gamma(1-\gamma))^n \end{aligned} \quad (45)$$

In fig. 9A we plot the values of the coefficients  $c$  as a function of  $\tilde{\gamma}$ . It is interesting to notice that both  $c_4$  and  $c_3$ , controlling the two leading orders, are null for the critical value  $\gamma = 1/2$ . This leaves the next leading order to  $c_2$  which, for this value of  $\gamma$ , is negative. As a result, we expect that the perturbation tends to lower the values of  $q_n$  at large  $n$  for  $\gamma \sim 1/2$ , and to raise it for  $\gamma \lesssim 1/2$ . This means that large-size extinction events are

made less probable by the presence of mutations around the critical value, which is consistent with the removal of the divergence observed in fig. 7C when mutations are present.

In fig. 9B,C,D we compare the above prediction for  $\Delta q_n$  with the value provided by numerical simulations for different values of  $\gamma$ , as indicated in each plot. We expect the theoretical prediction to be accurate for large values of  $n$ , and as long as the perturbation  $\Delta q_n$  remains small with respect to the unperturbed value  $q_n^0$ . However, for increasing values of  $n$  the perturbation grows faster than the unperturbed value, eventually invalidating this assumption. As a proxy for an accuracy upper limit we mark with a green dotted line the value of  $n$  at which the perturbation  $\Delta q$  has a magnitude equal to 10% of the unperturbed value  $q_n^0$ . Based on eqn. (45) this value scales as  $\sigma^{-4/5}$  for small  $\sigma$ 's.

## V. DISCUSSION

In this work we focused on the effects of a bottleneck on a B-cell population in the course of the affinity maturation process. Through a recursive relation that links the probability of bottleneck survival of a cell to the one of its daughter cells we were able to retrieve the dependence of a lineage extinction probability on its progenitor affinity. For lineages that go extinct we also evaluated the mean and variance of extinction time and progeny size, revealing a peak in extinction time corresponding to average affinity progenitors. Lineages stemming from these progenitors spawn in equilibrium between extinction and survival, and persist in this state until mutations drive the lineage either to survival or to extinction. Building on these results we then evaluated the survival probability for the full population as a function of Ag concentration and population size. We also included the effect of competition in an effective manner, using the deterministic model limit combined with a finite-size correction.

The bottleneck phenomenology was included in different maturation models [20, 21], which considered as optimal the maturation regime in which the B-cell population was subject to a strong enough selection force to grant good affinity enhancement, while at the same time not strong enough to cause population extinction. While the properties of the above models were numerically evaluated using stochastic simulations, we here present exact or approximate derivations for various quantities of interest through the use of recursive relations and probability generating functions. These techniques are often encountered in the context of branching processes, and similar approaches have been used in the study of AM [26, 27]. In our case we coupled the theoretical analysis with a more realistic model that included both the effect of mutations and selection for Ag binding and competition. We believe our approach provides a better understanding of what controls the lineage survival probability, while at the same time requiring less computa-

tional resources when compared to averaging over many stochastic simulations. This allows one to easily explore the effect of changing different model parameters on the survival probability and the lineage size.

In this last section of the paper we discuss the biological relevance of our results, together with some perspectives. We focus on four aspects: the role of population bottleneck as a limiting case in biologically observed maturation, the affinity of the initial B-cell population, how our results could be applied to explain the role of precursor frequency and affinity on the successful colonization of GCs, and finally how the theory might also offer insight in the case of multiple antigens.

### A. Population bottleneck in affinity maturation

The population bottleneck phenomenon, intended as a transitory state of low population with a high extinction risk, is a common feature of many maturation models [20, 21]. However to our knowledge this phenomenon lacks experimental confirmation. This might be due to two main factors: experimental limitations on one hand, and biologically relevant conditions on the other.

By nature the experimental observations of germinal centers tends to be destructive: in order to study the cellular composition of a germinal center the animal must often be sacrificed. This allows for the contemporary observation of multiple germinal centers at the same point in time, but not for the repeated observation of the same germinal center at different time points. As a result it is difficult to estimate for example whether chronic germinal center reactions feature long-lived germinal centers, or many short-lived germinal centers that appear and wane [8]. To our knowledge the only system capable of circumventing these difficulties is the one introduced by Firl and colleagues [23], who developed an intravital imaging system that allows for observation of germinal centers for an extended period of time. Interestingly, they report that in around 40% of their observations the germinal center failed to form, with the population of founder cells initially expanding until around 5 days after antigen administration, and then starting to wane. The authors observe that this might either be due to a competition with a non-imaged clone, or to a failure to establish a productive germinal center, which would be compatible with a "failed bottleneck" scenario. In short, while lacking explicit experimental confirmation, the bottleneck phenomenology is also not ruled out by experimental observations.

A second important consideration is that the bottleneck scenario might only represent a particular possible regime in maturation, which might be relatively rare under biologically relevant conditions. To illustrate this we consider the probability for a progenitor with energy  $\epsilon_{\text{progenitor}}$  to survive the bottleneck state, as a function

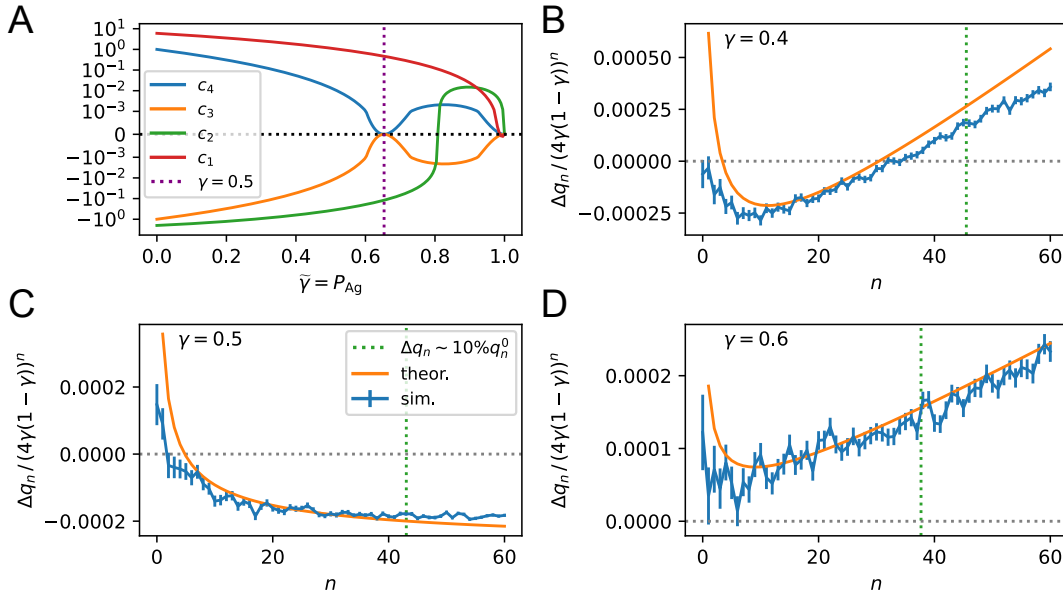


FIG. 9. **A**: values of the power expansion coefficients in eq. (43) as functions of the progenitor survival probability  $\tilde{\gamma} = P_{Ag}$ . Both  $c_4$  and  $c_3$  are null at the critical value  $\gamma = 1/2$ . **B,C,D**: perturbations to the probabilities  $q_n$  due to weak mutations (cf. eq. (36), with  $\sigma = 0.05$ ) for three different progenitor affinities, corresponding to  $\gamma = 0.4, 0.5, 0.6$  (B,C and D respectively). The perturbation  $\Delta q_n = q_n - q_n^0$  is evaluated as the difference between the probability of lineage extinction at progeny size  $n$  without ( $q_n^0$ ) and with ( $q_n$ ) mutations. These were evaluated by performing  $10^8$  numeric simulations for each of the two conditions. In the plots we report the value  $\Delta q_n (4\gamma(1-\gamma))^{-n}$  from simulations (blue) and compare it with the theoretical prediction (orange). The value of  $n$  at which  $|\Delta q_n| = 10\% |q_n^0|$  sets an upper limit for the validity of the theory (green).

of the amount of available antigen concentration and of the energy difference  $\epsilon_{\text{progenitor}} - \epsilon_{Ag}$  (cf. fig. 10A). This can be easily evaluated using the  $t \rightarrow \infty$  limit of recursion eq. (5). Here we only consider Ag-binding selection as the major source of selective pressure in the bottleneck and neglect the effect of competition for T-cell help. One can observe a steep transition between a region of almost certain extinction (red) to a region of possible survival (white and blue). The width of the region is essentially controlled by the probability of developing a beneficial mutation, and the magnitude of such mutations. As evident from the form of the selection survival probability eq. (2) operating a multiplicative change of antigen concentration by a factor  $\Delta C$  has the same effect as considering a change of binding energy of magnitude  $-\log \Delta C$ . One could therefore simply include the effect of concentration by rescaling the energy values. As a next step we consider the probability of survival for a full population, as a function of the initial population size  $N_i$  and the mean  $\mu_i$  and standard deviation  $\sigma_i$  of the initial binding energy distribution of cells. In fig. 10C we display the values at which the population survival probability  $p_{\text{surv}}$  is equal to 50%. These correspond to lines in the  $\mu_i / \sigma_i$  plane, with different colors corresponding to different initial population sizes  $N_i$ . The lines split the diagram in two regions, of high and low survival probability ( $p_{\text{surv}} \lesseqgtr 50\%$ ). To visualize the behavior of the system we consider three different cases, of high (case A), inter-

mediate (case B) or low (case C) survival probability. In fig. 10D for each of these conditions we display the evolution of the total population size for 100 different stochastic simulations. In the low-survival-probability region, corresponding to case C, in all of the simulations the population quickly goes extinct. Conversely, in the high-survival-probability region (case A), the population only undergoes a partial reduction in size, with all of the simulations successfully surpassing the bottleneck. In the intermediate region (case C) we observe a consistent reduction in size for all of the simulations, with both successes and failures in overcoming the bottleneck. This intermediate region corresponds indeed to the bottleneck phenomenology that we are interested in studying in this paper, but biologically might only represent a subset of all the possible maturation regimes. In particular one might expect maturation to normally occur in the high-survival-probability regime in most of the biologically relevant cases, especially for pathogens that are similar to ones already encountered in the past for which high-affinity precursor might already be available in the initial population. The low-survival-probability regime conversely might never occur. In this case in fact no precursor with sufficient affinity is expected to be present in the initial population, and GCs might fail to form altogether. Finally the intermediate regime, in which the bottleneck phenomenology is observed, might be a limit case and only occur for “hard” antigens, when only few precursor of acceptable affinity exist in

the initial population. These limit cases might be of particular interest when developing novel vaccination strategies, for example against mutable pathogens such as HIV, as will be discussed in section V C.

### B. Affinity of the initial population

As shown in the previous section, the bottleneck phenomenology occurs only for values of the parameters that place the system in the boundary between the high and low survival probability region. This is fundamentally controlled by the affinity of the initial population: if enough high-affinity precursors are present in the initial population then the population survival is almost certain, while if all cells in the population have low affinity then survival is rare event. To study this phenomenology in our simulations we chose the size and affinity distribution of the initial population so as to be on this boundary region, while at the same time being compatible with experiments. This was done as follows. The initial population size was set equal  $N_i = 2500$ . This is in agreement with [5] which reports around 3000 cells per germinal center. For simplicity we consider the binding energy of each cell to be independently extracted from a normal distribution with mean  $\mu_i$  and standard deviation  $\sigma_i$ . This simplification overestimates the heterogeneity of the initial population binding energy, which in reality is usually composed of  $50 \sim 200$  different clones [28]. The value of  $\sigma_i = 1.5$  was chosen in agreement with [13], in which a fit of experimental single-cell affinity measurements of naive responder cells results in an estimated  $\sigma_i \sim 1.66$ . Given these conditions the value of  $\mu_i$  (together with antigen concentration  $C$ ) controls the probability of surviving the bottleneck (cf. fig. 10C). For most of the simulation we choose a value  $\mu_i = 4$ , so that only around 2% of cells from the initial population have affinity higher than the Ag-binding selection threshold for values of  $C \sim 5$ . This makes so that for standard parameter values the model is close to the phase threshold line (i.e. condition B in fig. 10C and D), which allows us to visualize and study the bottleneck phenomenology.

Under these assumptions the initial population will be composed in large part by low-affinity cells. In general, while antigen affinity does indeed drive the recruitment of cells in germinal centers, the presence of cells with undetectable affinity has been observed in different experiments [16]. For example in the experiments reported in [29] only a minor fraction ( $\sim 30\%$ ) of all the antibody-secreting cells elicited by immunization was showing detectable affinity for the administered antigen ( $K_d < 500$  nM). Moreover, as experimentally shown in [30], in the absence of competition even cells with very low affinity for the antigen ( $K_d \sim 8$   $\mu$ M) can colonize GCs. The role of low-affinity cells in maturation has not yet been experimentally elucidated. Their presence

has either been attributed to some form of non-specific activation (e.g. bystander activation [31, 32]), or to the specific binding of some unknown “dark antigen” [16], or again to specific binding of the administered antigen but with a very low (undetectable) affinity. In either case their number usually decreases during maturation, in favor of higher affinity cells.

### C. Effect of affinity and precursor frequency on germinal center colonization

The theory we introduced in this paper offers a framework to quantify the probability for a cell lineage to be able to survive the population bottleneck, successfully colonize a GC and undergo maturation. This has strong connections with the field of vaccine design. In this field the development of new and effective vaccination techniques against pathogens such as HIV requires the stimulation of germline clones that are oftentimes rare or have low affinity for the target. To design an effective vaccine it is indeed important to understand the effect of germline affinity and precursor frequency on the successful recruitment of clones in germinal centers. Recent experiments have shown that at physiological levels of affinity and frequency, successful recruitment depends in a non-trivial way on both of these variables [15, 16, 33].

In this context, our theory could offer a quantitative prediction for the probability of a clonal family to survive selection, as a function of the initial abundance of cells and of their affinity. To illustrate this we consider a group of  $N$  identical cells with initial energy  $\epsilon_{\text{clone}}$ . Using our theory we can evaluate the probability that the offspring of at least one cell will survive the population bottleneck. This probability is displayed in fig. 10B. In agreement with what observed in experiments [15, 16], successful survival depends on both  $N$  and  $\epsilon_{\text{clone}}$ , with survival probability being increased by both higher affinity and higher precursor frequency. A quantitative understanding of this trade-off might guide decisions in the development of vaccine antigens and improve experiment design.

### D. Case of multiple antigens

It would be important to extend the present work to the case of multiple antigens. Contrary to the case of simple antigens, where the response is usually focused on a single dominating epitope, in the case of complex or multiple antigens there might be competition between the binding of different epitopes. Understanding how maturation plays out in the presence of multiple antigens is currently an open issue, whose solution could help the development of vaccination strategies against muta-



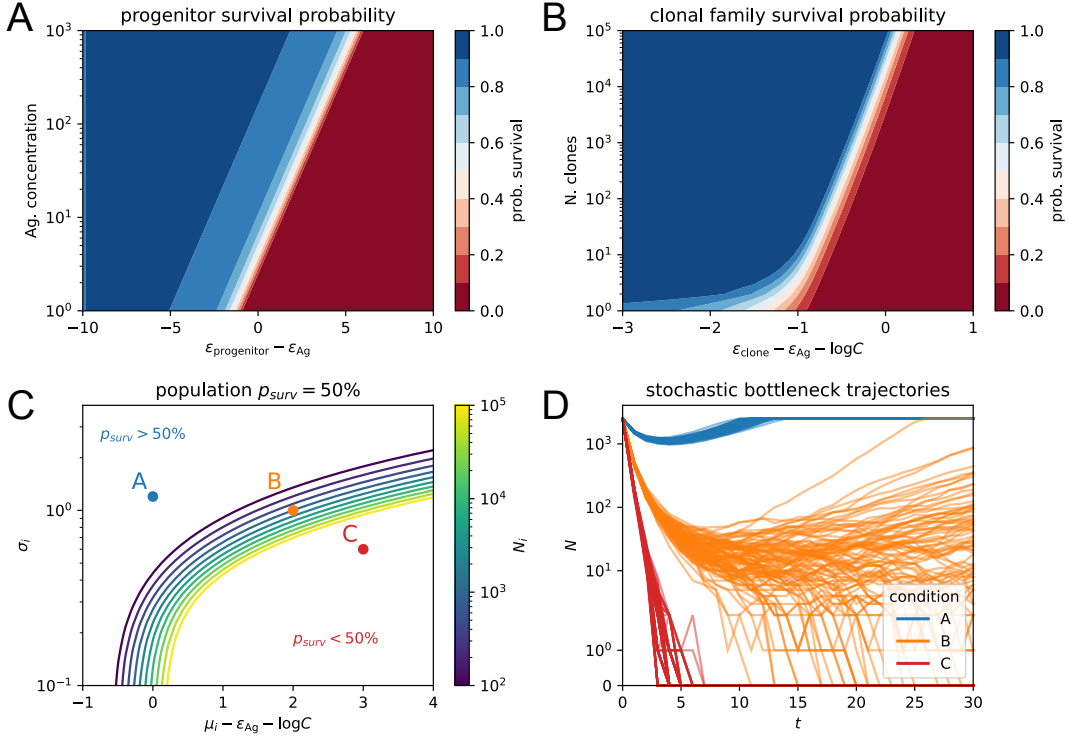


FIG. 10. **A:** bottleneck survival probability of a lineage stemming from a single progenitor cell with energy  $\epsilon_{\text{progenitor}}$ , as a function of antigen concentration. **B:** probability of bottleneck survival for an initial population composed of  $N$  identical cells with binding energy  $\epsilon_{\text{clone}}$ . **C:** population survival probability as a function of its initial size  $N_i$  and the mean  $\mu_i$  and standard deviation  $\sigma_i$  of the initial binding energy distribution. Colored lines indicate the region of the phase diagram in which the probability of survival is equal to 50%. Different colors distinguish different initial population sizes according to the colorbar on the right. Each line divides the plane in two regions, corresponding to high ( $p_{\text{surv}} > 50\%$ ) and low ( $p_{\text{surv}} < 50\%$ ) survival probability. We consider three conditions that belong to either one of these two regions (condition A and C), or to the separation line (condition B) to be visualized in the next panel. **D:** evolution of total population size  $N$  for stochastic model simulations realized under three different conditions as indicated in panel C. For each condition we display 100 different stochastic simulation trajectories, with time  $t$  measured in evolution rounds. The simulation were performed at an initial population size  $N_i = 2500$  and antigen concentration  $C = 10$ , while the mean and standard deviation ( $\mu_i, \sigma_i$ ) for the initial binding energy distribution were varied according to the condition (A: (2.3, 1.2), B: (4.3, 1.), C: (5.3, 0.6)). Results displayed in panels A, B and C do not take into account competitive selection, while in panel D competitive selection is included in simulations.

ble pathogens such as HIV [8, 34, 35].

As a natural extension to our theory one could consider that a single cell might possess different affinities for each antigen mutant. These affinities are potentially correlated, depending on the similarity between the different mutants, and so is the effect of mutations. Taking this into account one could then define mutation and selection probabilities in the presence of multiple Ag mutants, and use a similar approach to the one introduced in this paper to evaluate the lineage and population survival probability.

Even without such extension however we believe that the single-antigen framework might still be instructive in the study of multiple antigen maturation. To demonstrate this one can extend the model to the case of two antigens, in a similar fashion to what was done in [21]. In this extension each cell possesses two binding energies  $\epsilon_1$  and  $\epsilon_2$ , encoding its affinity for the two antigens. The val-

ues of these energies are initially independently extracted from the same normal distribution of "naive" binding energies. The system then evolves similarly to the single-antigen case, with two differences. The first one is that affinity-affecting mutations change the binding energy of the cell for both antigens with two random contributions  $\Delta\epsilon_1$  and  $\Delta\epsilon_2$ , independently extracted from the usual mutation kernel (cf. eq. (1)). Notice that as a result of this, mutations that increase the affinity for only one of the two antigens are much more common than mutations that increase the affinity for both. The second difference is in the way selection is performed. Following the "meet-all" scenario introduced in [21] we extend the definition

of the selection survival probabilities as:

$$P_{\text{Ag}}(\epsilon_1, \epsilon_2) = \frac{C_1 e^{-\epsilon_1} + C_2 e^{-\epsilon_2}}{C_1 e^{-\epsilon_1} + C_2 e^{-\epsilon_2} + e^{-\epsilon_{\text{Ag}}}} \quad (46)$$

$$P_{\text{T}}(\epsilon_1, \epsilon_2) = \frac{C_1 e^{-\epsilon_1} + C_2 e^{-\epsilon_2}}{C_1 e^{-\epsilon_1} + C_2 e^{-\epsilon_2} + e^{-\bar{\epsilon}}}, \quad (47)$$

with  $e^{-\bar{\epsilon}} = \langle e^{-\epsilon_1} + e^{-\epsilon_2} \rangle_{\text{pop}}$

Where  $C_1$  and  $C_2$  represent the concentrations of the two antigens. As a result of these definitions a cell will survive selection if it is able to bind at least one of the two antigens with good affinity.

The combination of these two effects, namely the fact that beneficial mutation will often affect only one of the two components, and selection will allow for survival of cells that have high affinity for at least one of the two antigens, makes so that evolution will split the population in two subsets: cells with high affinity for one antigen or cells with high affinity for the other. To visualize this we simulated the evolution of the system at antigen concentration  $C_1 = C_2 = 5$ . We performed 500 different simulations. In fig. 11 we plot the evolution of the average binding energy distribution  $P_t(\epsilon_1, \epsilon_2)$  of the population, and the two marginal distributions  $P_t(\epsilon_1)$ ,  $P_t(\epsilon_2)$  at four different times  $t = 0, 10, 50, 100$ . The marginal distributions show how the population divides in these two subsets. These subsets evolve independently, each one improving its affinity for their “target” antigen over time, while the binding energy for the other antigen gradually decreases due to deleterious mutations not being selected against. The evolution of each of these subsets is therefore well captured by the single-antigen framework along the corresponding axis.

We point out that this split in the population is a consequence of the two assumptions we made, on the way mutations and selection operate on the two energy component. Evolution can proceed differently if for example the effect of mutations on the two components is strongly correlated, or if selection requires cells to have high affinity for all of the antigens at the same time.<sup>2</sup> If high affinity for both antigens is required then the population will not split in two subsets, but rather evolve along the diagonal of the affinity space  $\epsilon_1 = \epsilon_2$ . One could therefore simply study the dynamics along a single component, for which the single antigen case would still offer a useful framework.

## E. Model limitations and perspectives

Compared to real AM, our model is obviously simplified in many aspects, for example we do not impose

---

<sup>2</sup> The latter case seems less likely given the presence of low-affinity responders in immunization experiments [29, 36].

an affinity ceiling and beneficial mutations can accumulate indefinitely. However we believe this approximation to be reasonable since for the bottleneck scenario we consider low-affinity cells that could potentially undergo many affinity-improving mutations. Moreover we consider Ag concentration to be constant, while in reality Ag is subject to natural decay and consumption by B-cells. We believe this approximation to be acceptable when studying bottleneck survival, that in most cases is resolved in few evolution rounds. In our analysis we focused on the evolution of the population of B-cells inside a single GC. In the course of a response however many GCs form inside the body.<sup>3</sup> By averaging over the distribution of energies in the initial population our theory indeed quantifies the average bottleneck survival probability of GCs in this ensemble, in the hypothesis that the affinity of cells each initial population is independent. If instead affinities are correlated or if cells are able to migrate between GCs then the average would be harder to compute. Unfortunately the lack of precise experimental quantification of these processes forbids any meaningful modeling so far.

Finally, our model does not account for the fact that GC selection might be highly permissive. Indeed, low affinity cells have been shown to be able to reside in GCs for an extended periods of time, and permissiveness might especially characterize maturation against complex antigens [13, 38]. This effect could be accounted for in our model by adding a small affinity-independent probability of surviving selection. This might have a beneficial impact on maturation, since it could allow low-affinity cells to reside longer in GCs and it might provide them a higher chance of developing beneficial mutations.

**Acknowledgments:** We are deeply grateful to Jean Baudry, Arup Chakraborty and Klaus Eyer for many useful discussions and interactions.

## Appendix A: Model parameters choice

The values of the parameters are reported in table I, and were chosen based on existing literature.

Mature GCs have a B-cells population consisting of a few thousands cell [28, 39, 40]. We therefore set the initial and maximum size of the population equal to  $N_i = N_{\text{max}} = 2500$ . This is in agreement with [5] which reports around 3000 cells per GC. However we stress that GCs are heterogeneous in size [37]. Similarly to [13, 21] we consider the duration of a turn of evolution to be  $T_{\text{round}} \sim 12\text{h}$ , which is consistent with timing of cell migration [4, 41]. Time in our model will be rescaled by this

---

<sup>3</sup> While their number is not known with accuracy, it could range from many tens to few hundreds since spleen sections revealed around 20-50 GCs in mice [37].

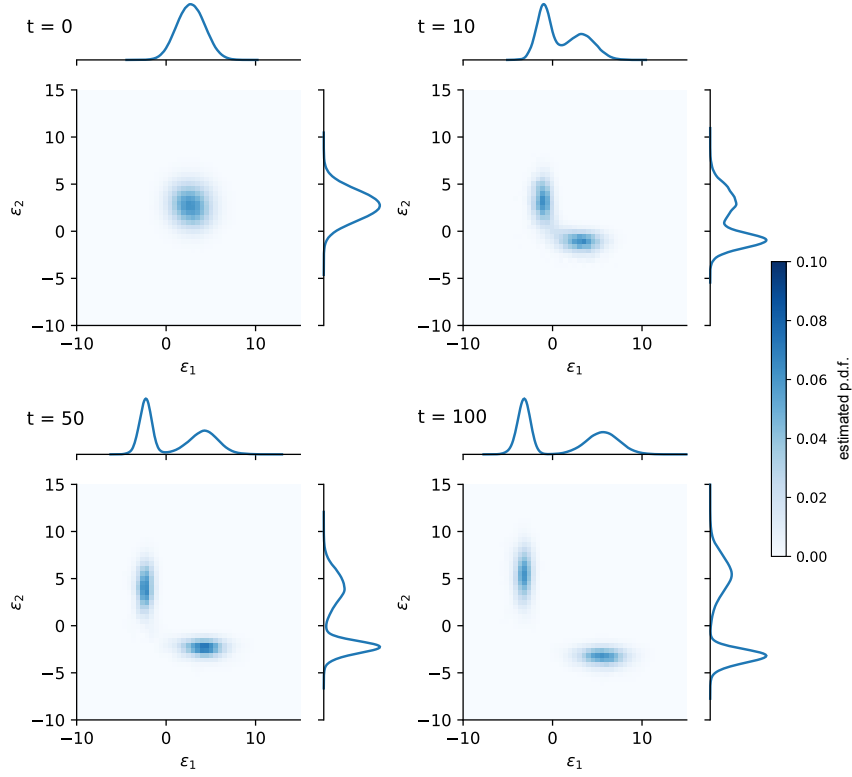


FIG. 11. Joint probability distribution  $P_t(\epsilon_1, \epsilon_2)$  of the binding energies  $(\epsilon_1, \epsilon_2)$  of the B-cell population in the two-antigens case described in the text. The four different plots correspond to four different times ( $t = 0, 10, 50, 100$ ), measured in evolution rounds. The probability distribution was estimated by averaging over 500 different stochastic simulations. For each plot we also display the marginal probabilities  $P_t(\epsilon_1)$  and  $P_t(\epsilon_2)$ . Simulations were performed at a value of antigen concentration  $C_1 = C_2 = 5$ , and considering an initial average binding energy  $\mu_i = 3$  for both antigens.

parameter	value	description
$\mu_i, \sigma_i$	4, 1.5	mean and variance of the population initial binding energy distribution
$N_i$	2500	initial population size
$N_{\max}$	2500	maximum carrying capacity
$p_{\text{sil}}, p_{\text{let}}, p_{\text{aa}}$	0.75, 0.15, 0.1	probabilities of silent, lethal, affinity-affecting mutations
$\mu_M, \sigma_M$	0.3, 0.3	mean and variance of distribution of affinity-affecting mutations
$\epsilon_{\text{Ag}}$	0	Ag-binding selection threshold energy
$C$	see figure captions	Ag concentration
$p_{\text{diff}}$	0.1	differentiation probability

TABLE I. standard values of model parameters. Unless otherwise specified these are the values used in simulations. The choice of their value is discussed in appendix A

standard quantity, so that the variable  $t$  has no dimension. Similarly, also the binding energy  $\epsilon$  is dimensionless, expressed in standard units of  $k_B T$ . For simplicity following experiments that indicate a cell-cycle time of 12h or longer [42] we consider a single division per round. Other experiments indicate an average of two cell division in the DZ [43]. We point out that, at the expense of simplicity, our theory can also be extended to account for an higher number of cell divisions.

In [20, 22] mutations occur at a rate of 0.5 per sequence per division, and are silent, lethal or affinity affecting with probabilities of respectively 0.5, 0.3, 0.2. This fixes

our effective mutation probabilities to  $p_{\text{sil}} = 0.75, p_{\text{let}} = 0.15, p_{\text{aa}} = 0.1$ . To reproduce the fact that most of the mutations are deleterious we pick for simplicity  $\mu_M = \sigma_M$ . This fixes the amount of beneficial mutations to  $\sim 16\%$ , which is somewhat higher but still compatible with other models [13, 20, 21] in which this fraction is set to 5%. Moreover we set  $\mu_M = 0.3$  so as to set the mean effect of beneficial mutations to  $\langle \Delta \epsilon \rangle_{\text{beneficial}} \sim -0.15$ . This value is slightly smaller than  $\langle \Delta \epsilon \rangle_{\text{beneficial}} \sim -0.53$  used in [13], but this is compensated by the higher rate of beneficial mutations in our model. As described in the main text, the binding energy distribution of the initial population

is set to a Gaussian with standard deviation  $\sigma_i = 1.5$ , which is compatible with experimental data [13]. Since evolution of the population is invariant for shifts of the energy space we set  $\epsilon_{\text{Ag}} = 0$ . Under this choice of gauge the zero in the energy space is the threshold energy for Ag-binding selection. Moreover we pick  $\mu_i = 4$  so that the difference  $\mu_i - \epsilon_{\text{Ag}} - \log C \sim 2\sigma_i$  for the values of Ag concentrations considered in this work ( $C \sim 5$ ) and on average only around 2% of cells from the initial population meet this threshold. For simplicity we independently extract the energy of each cell in the initial population from this initial distribution. By doing so we might overestimate the diversity of the initial population. In fact, experiments probing the clonal composition of GCs estimated that early GCs contain around 50 to 200 different clonal families [28]. A more realistic initiation of our GCs would require us to extract the energies of around a hundred different founder cells, and let them duplicate without mutating up to to the full GC capacity. This would generate a less homogeneous initial population than the one we consider in our simulation, but would not otherwise strongly impact our results. Lastly, as in [13, 21], the probability of differentiation is set to  $p_{\text{diff}} = 0.1$ .

### Appendix B: Estimation of $\bar{\epsilon}$ evolution

Including the effect of competition in our evaluation of the population survival probability requires us to estimate the evolution of  $\bar{\epsilon}$ , defined in eq. (3). To obtain a tractable approximation, we consider the deterministic limit of big population size. In this limit the population binding energy can be approximated with a continuous distribution, and the state of the system is completely determined by the density function  $\rho_t(\epsilon)$ . This function represents the density of cells having energy  $\epsilon$  at evolution round  $t$ , so that its integral is equal to the size of the population, and its normalized version is the population binding energy distribution. Evolution is expressed in terms of operators acting on this function. In particular:

1. Cell duplication corresponds simply to doubling in size:

$$\mathbf{A}[\rho](\epsilon) = 2 \times \rho(\epsilon) \quad (\text{B1})$$

2. Mutations are represented as the convolution with the mutation kernel  $K(\Delta\epsilon)$  defined in eq. (1). Notice that the kernel  $K$  is not normalized, to account for the contribution of lethal mutations. It acts on the distribution as:

$$\mathbf{M}[\rho](\epsilon) = \int d\Delta\epsilon \rho(\epsilon - \Delta\epsilon) K(\Delta\epsilon) \quad (\text{B2})$$

3. Ag-binding selection is implemented by in the product of the population function with the survival probability eq. (2):

$$\mathbf{S}_{\text{Ag}}[\rho](\epsilon) = P_{\text{Ag}}(\epsilon) \rho(\epsilon) \quad (\text{B3})$$

4. Similarly, T-cell help selection is given by the product with the survival probability eq. (3):

$$\mathbf{S}_{\text{T}}[\rho](\epsilon) = P_{\text{T}}(\epsilon, \bar{\epsilon}) \rho(\epsilon), \quad \text{with } e^{-\bar{\epsilon}} = \frac{1}{N} \int d\epsilon e^{-\epsilon} \rho(\epsilon) \quad (\text{B4})$$

Where  $N = \int d\epsilon \rho(\epsilon)$  is the current population size.

5. Differentiation consists simply in a product involving the differentiation probability:

$$\mathbf{D}[\rho](\epsilon) = (1 - p_{\text{diff}}) \rho(\epsilon) \quad (\text{B5})$$

6. Finally, the carrying capacity constraint is implemented again by a product and is operated only if the size of the population exceeds the maximum limit:

$$\mathbf{C}[\rho](\epsilon) = \min\{1, N_{\text{max}}/N\} \rho(\epsilon) \quad (\text{B6})$$

Where again  $N = \int d\epsilon \rho(\epsilon)$  is the current population size.

From these definitions the evolution of the population density function  $\rho_t(\epsilon)$  can be expressed as:

$$\rho_{t+1} = \mathbf{C} \mathbf{D} \mathbf{S}_{\text{T}} \mathbf{S}_{\text{Ag}} \mathbf{M} \mathbf{A} \rho_t \quad (\text{B7})$$

Combining with the definition for  $\bar{\epsilon}$  eq. (B4) provides a way for us to estimate the average evolution of  $\bar{\epsilon}_t$ .<sup>4</sup>

### Appendix C: Derivation and numerical evaluation of recursion equation

The core of our theory rests on recursion eq. (5), which allows one to evaluate the probability of extinction of the progeny of a cell from the extinction probabilities of its daughter cells. In this appendix we present an intuitive derivation of this equation, and explain how it was numerically evaluated.

As described in the main text, eq. (5) operates a recursion on the function  $d_t(\epsilon)$ , defined as the probability that all the progeny of a given progenitor cell with affinity  $\epsilon$  will go extinct by evolution round  $t$ . The recursion consists in expressing this probability as a function of the progeny extinction probability for the two daughter cells of the progenitor in the remaining  $t - 1$  rounds. These daughter cells might develop a mutation which introduces a change in binding energy of magnitude  $\Delta\epsilon$ , and the recursion will therefore contain the term  $d_{t-1}(\epsilon + \Delta\epsilon)$ . The derivation proceeds in the following manner, schematized in eq. (C1). The

<sup>4</sup> Notice that from the order of the operators in the evolution round it follows that  $\epsilon_t$  must be evaluated using eq. (B4) not directly on  $\rho_t$  but rather on  $\mathbf{S}_{\text{Ag}} \mathbf{M} \mathbf{A} \rho_t$

probability that all of the progeny of the progenitor cell with energy  $\epsilon$  goes extinct by evolution round  $t$  is equal to the probability that both of its daughter cell progenies will go extinct before this round. These probabilities can be expressed as one minus their complement: the probability that extinction does not occur by round  $t$ . In turn this is equal to the probability that the daughter cell will mutate, survive selection of the first round, and some of its progeny will also survive the remaining  $t - 1$

rounds. Expressing this latter survival probability in terms of its complement completes the recursive relation. Notice that while the recursion is back in time ( $t \rightarrow t - 1$ ), the term  $d_{t-1}$  refers to the survival probability of the daughter cell, and in this sense the recursion is forward in time. In fact  $t$  refers in this case to the number of turns left to survive, and decreases of one unit at every jump to the next generation.

$$\begin{aligned}
 d_t(\epsilon) &= P(\text{progeny extinct by round } t) \\
 &= [P(\text{daughter cell's progeny extinct by round } t)]^2 \\
 &= [1 - P(\text{daughter cell's progeny not extinct by round } t)]^2 \\
 &= \left[ 1 - \int d\Delta\epsilon P(\text{mutation } \Delta\epsilon) P(\text{survive round}) P(\text{progeny survives next } t - 1 \text{ rounds}) \right]^2 \\
 &= \left[ 1 - \int d\Delta\epsilon P(\text{mutation } \Delta\epsilon) P(\text{survive round}) (1 - d_{t-1}(\epsilon + \Delta\epsilon)) \right]^2
 \end{aligned} \tag{C1}$$

Numerical evaluation of eq. (5) was performed according to the following scheme. The fact that cells can accumulate mutations makes so that the recursion requires evaluating an integral over all the possible values of the energy. As a consequence in order to evaluate  $d_t(\epsilon)$  one needs to know the value of  $d_{t-1}(\epsilon)$  at the previous time interval for any value of  $\epsilon$ . In practice, when evaluating the function at a given final time  $T$  one can define a matrix  $d[t, i]$ , where the first index  $t = 1, \dots, T$  runs over all possible values of times, and the index  $i = 0, \dots, I$  operates a discretization of the energy space  $\epsilon_0, \dots, \epsilon_I$ , in such a way that  $d[t, i] = d_t(\epsilon_i)$ . We indicate with

$\delta\epsilon = \epsilon_{i+1} - \epsilon_i$  the value of the discretization step. The first row of the matrix is given simply by the discretized form of eq. (4) in the main text:

$$d[1, i] = \left[ 1 - \sum_j \delta\epsilon K[j] P_S[i + j] (1 - p_{\text{diff}}) \right]^2 \tag{C2}$$

Where  $K[i] = K(i \delta\epsilon)$  and  $P_S[i] = P_S(\epsilon_i)$  are the discretized forms of the mutation kernel and selection survival probability. From here rows of the matrix can be iteratively populated using the discretized form of eq. (5):

$$d[t, i] = \left[ 1 - \sum_j \delta\epsilon K[j] P_S[i + j] (1 - p_{\text{diff}}) (1 - d[t - 1, i + j]) \right]^2 \tag{C3}$$

In practice the evaluation of the extinction probability at time  $T$  requires of the order of  $T$  convolutions between arrays having a size  $(\epsilon_I - \epsilon_0)/\delta\epsilon$ . Many efficient algorithms exist to perform fast numerical convolution. In our implementation for example, written in Python version 3.9, the evaluation of  $d_t$  for  $t = 150$  and discretizing the binding energy interval  $[-50, 50]$  at a discretization step  $\delta\epsilon = 0.01$  requires around 0.15 seconds on a laptop. The more general case of a time-varying survival probability, described in eqs. (12) and (13) in the main text, also presents a similar complexity. In this case one is in-

terested in evaluating the probability  $d_{0,T}(\epsilon)$  that a cell having energy  $\epsilon$  at time 0 will have its progeny extinct by evolution round  $T$ . In order to do so one defines a matrix  $d[t, i] = d_{t,T}(\epsilon_i)$ , with  $t = 0, \dots, T - 1$ . The row of the matrix for  $t = T - 1$  can be populated from the discretized form of eq. (12) in the main text. At this point using the discretized form of eq. (13) one can populate row  $t - 1$  from row  $t$ , until reaching the desired value  $t = 0$ . As in the previous case, this requires of the order of  $T$  convolutions of arrays having the size of the discretized binding energy space, and requires a similar computational time.

- [1] G. D. Victora and M. C. Nussenzweig, Germinal centers, *Annual review of immunology* **30**, 429 (2012).
- [2] N. S. De Silva and U. Klein, Dynamics of b cells in germinal centres, *Nature reviews immunology* **15**, 137 (2015).
- [3] O. Bannard and J. G. Cyster, Germinal centers: programmed for affinity maturation and antibody diversification, *Current Opinion in Immunology* **45**, 21 (2017).
- [4] L. Mesin, J. Ersching, and G. D. Victora, Germinal center B cell dynamics, *Immunity* **45**, 471 (2016).
- [5] H. N. Eisen, Affinity enhancement of antibodies: how low-affinity antibodies produced early in immune responses are followed by high-affinity antibodies later and in memory B-cell responses, *Cancer immunology research* **2**, 381 (2014).
- [6] G. D. Victora and L. Mesin, Clonal and cellular dynamics in germinal centers, *Current Opinion in Immunology* **28**, 90 (2014).
- [7] M. J. Shlomchik, W. Luo, and F. Weisel, Linking signaling and selection in the germinal center, *Immunological reviews* **288**, 49 (2019).
- [8] G. D. Victora and H. Mouquet, What are the primary limitations in b-cell affinity maturation, and how much affinity maturation can we drive with vaccination? lessons from the antibody response to hiv-1, *Cold Spring Harbor perspectives in biology* **10**, a029389 (2018).
- [9] S. H. Kleinstein, Y. Louzoun, and M. J. Shlomchik, Estimating hypermutation rates from clonal tree data, *The Journal of Immunology* **171**, 4639 (2003).
- [10] D. McKean, K. Huppi, M. Bell, L. Staudt, W. Gerhard, and M. Weigert, Generation of antibody diversity in the immune response of BALB/c mice to influenza virus hemagglutinin., *Proceedings of the National Academy of Sciences* **81**, 3180 (1984).
- [11] C. Berek and C. Milstein, Mutation drift and repertoire shift in the maturation of the immune response, *Immunological reviews* **96**, 23 (1987).
- [12] S. J. Rhodes, G. M. Knight, D. E. Kirschner, R. G. White, and T. G. Evans, Dose finding for new vaccines: the role for immunostimulation/immunodynamic modelling, *Journal of theoretical biology* **465**, 51 (2019).
- [13] M. Molari, K. Eyer, J. Baudry, S. Cocco, and R. Monasson, Quantitative modeling of the effect of antigen dosage on B-cell affinity distributions in maturing germinal centers, *eLife* **9**, e55678 (2020).
- [14] M. Kang, T. J. Eisen, E. A. Eisen, A. K. Chakraborty, and H. N. Eisen, Affinity inequality among serum antibodies that originate in lymphoid germinal centers, *PLOS ONE* **10**, e0139222 (2015).
- [15] R. K. Abbott, J. H. Lee, S. Menis, P. Skog, M. Rossi, T. Ota, D. W. Kulp, D. Bhullar, O. Kalyuzhniy, C. Havenar-Daughton, *et al.*, Precursor frequency and affinity determine b cell competitive fitness in germinal centers, tested with germline-targeting hiv vaccine immunogens, *Immunity* **48**, 133 (2018).
- [16] R. K. Abbott and S. Crotty, Factors in b cell competition and immunodominance, *Immunological reviews* **296**, 120 (2020).
- [17] A. K. Chakraborty, A perspective on the role of computational models in immunology, *Annual review of immunology* **35**, 403 (2017).
- [18] L. Buchauer and H. Wardemann, Calculating germinal centre reactions, *Current Opinion in Systems Biology* **18**, 1 (2019).
- [19] F. Horns, C. Vollmers, C. L. Dekker, and S. R. Quake, Signatures of selection in the human antibody repertoire: Selective sweeps, competing subclones, and neutral drift, *Proceedings of the National Academy of Sciences* **116**, 1261 (2019).
- [20] J. Zhang and E. I. Shakhnovich, Optimality of mutation and selection in germinal centers, *PLoS computational biology* **6**, e1000800 (2010).
- [21] S. Wang, J. Mata-Fink, B. Kriegsman, M. Hanson, D. J. Irvine, H. N. Eisen, D. R. Burton, K. D. Wittrup, M. Kardar, and A. K. Chakraborty, Manipulating the selection forces during affinity maturation to generate cross-reactive hiv antibodies, *Cell* **160**, 785 (2015).
- [22] S. Wang, Optimal sequential immunization can focus antibody responses against diversity loss and distraction, *PLoS computational biology* **13**, e1005336 (2017).
- [23] D. J. Firl, S. E. Degn, T. Padera, and M. C. Carroll, Capturing change in clonal composition amongst single mouse germinal centers, *Elife* **7**, e33051 (2018).
- [24] H. W. Watson and F. Galton, On the probability of the extinction of families, *The Journal of the Anthropological Institute of Great Britain and Ireland* **4**, 138 (1875).
- [25] P. Flajolet and R. Sedgewick, *Analytic combinatorics* (Cambridge University Press, 2009).
- [26] I. Balelli, V. Milišić, and G. Wainrib, Random walks on binary strings applied to the somatic hypermutation of b-cells, *Mathematical biosciences* **300**, 168 (2018).
- [27] I. Balelli, V. Milišić, and G. Wainrib, Multi-type galton-watson processes with affinity-dependent selection applied to antibody affinity maturation, *Bulletin of mathematical biology* **81**, 830 (2019).
- [28] J. M. Tas, L. Mesin, G. Pasqual, S. Targ, J. T. Jacobsen, Y. M. Mano, C. S. Chen, J.-C. Weill, C.-A. Reynaud, E. P. Browne, *et al.*, Visualizing antibody affinity maturation in germinal centers, *Science* **351**, 1048 (2016).
- [29] K. Eyer, C. Castrillon, G. Chenon, J. Bibette, P. Bruhns, A. D. Griffiths, and J. Baudry, The quantitative assessment of the secreted igg repertoire after recall to evaluate the quality of immunizations, *The Journal of Immunology* **205**, 1176 (2020).
- [30] T. A. Schwickert, G. D. Victora, D. R. Fooksman, A. O. Kamphorst, M. R. Mugnier, A. D. Gitlin, M. L. Dustin, and M. C. Nussenzweig, A dynamic T cell-limited checkpoint regulates affinity-dependent B cell entry into the germinal center, *Journal of Experimental Medicine* **208**, 1243 (2011).
- [31] T. Inoue, I. Moran, R. Shinnakasu, T. G. Phan, and T. Kurosaki, Generation of memory B cells and their reactivation, *Immunological reviews* **283**, 138 (2018).
- [32] F. Horns, C. L. Dekker, and S. R. Quake, Memory b cell activation, broad anti-influenza antibodies, and bystander activation revealed by single-cell transcriptomics, *Cell reports* **30**, 905 (2020).
- [33] C. Havenar-Daughton, R. K. Abbott, W. R. Schief, and S. Crotty, When designing vaccines, consider the starting material: the human b cell repertoire, *Current opinion in immunology* **53**, 209 (2018).
- [34] A. K. Chakraborty and J. P. Barton, Rational design of vaccine targets and strategies for HIV: a crossroad

- of statistical physics, biology, and medicine, *Reports on Progress in Physics* **80**, 032601 (2017).
- [35] P. A. Robert, A. L. Marschall, and M. Meyer-Hermann, Induction of broadly neutralizing antibodies in germinal centre simulations, *Current opinion in biotechnology* **51**, 137 (2018).
- [36] M. Kuraoka, A. G. Schmidt, T. Nojima, F. Feng, A. Watanabe, D. Kitamura, S. C. Harrison, T. B. Kepler, and G. Kelsoe, Complex antigens drive permissive clonal selection in germinal centers, *Immunity* **44**, 542 (2016).
- [37] N. Wittenbrink, A. Klein, A. A. Weiser, J. Schuchhardt, and M. Or-Guil, Is there a typical germinal center? a large-scale immunohistological study on the cellular composition of germinal centers during the hapten-carrier-driven primary immune response in mice, *The Journal of Immunology* **187**, 6185 (2011).
- [38] J. Finney, C.-H. Yeh, G. Kelsoe, and M. Kuraoka, Germinal center responses to complex antigens, *Immunological reviews* **284**, 42 (2018).
- [39] J. Jacob, J. Przylepa, C. Miller, and G. Kelsoe, In situ studies of the primary immune response to (4-hydroxy-3-nitrophenyl) acetyl. III. the kinetics of V region mutation and selection in germinal center B cells., *Journal of Experimental Medicine* **178**, 1293 (1993).
- [40] M. McHeyzer-Williams, M. McLean, P. Lalor, and G. Nossal, Antigen-driven B cell differentiation in vivo., *Journal of Experimental Medicine* **178**, 295 (1993).
- [41] G. D. Victora, T. A. Schwickert, D. R. Fooksman, A. O. Kamphorst, M. Meyer-Hermann, M. L. Dustin, and M. C. Nussenzweig, Germinal center dynamics revealed by multiphoton microscopy with a photoactivatable fluorescent reporter, *Cell* **143**, 592 (2010).
- [42] C. D. Allen, T. Okada, H. L. Tang, and J. G. Cyster, Imaging of germinal center selection events during affinity maturation, *Science* **315**, 528 (2007).
- [43] A. D. Gitlin, Z. Shulman, and M. C. Nussenzweig, Clonal selection in the germinal centre by regulated proliferation and hypermutation, *Nature* **509**, 637 (2014).