

1 **Random walks on binary strings applied to the**
2 **somatic hypermutation of B-cells**

3 Irene Balelli · Vuk Milišić · Gilles Wainrib

4
5 Received: date / Accepted: date

6 **Abstract** Within the germinal center in follicles, B-cells proliferate, mutate
7 and differentiate, while being submitted to a powerful selection: a micro-
8 evolutionary mechanism at the heart of adaptive immunity. A new foreign
9 pathogen is confronted to our immune system, the mutation mechanism that
10 allows B-cells to adapt to it is called *somatic hypermutation*: a programmed
11 process of mutation affecting B-cell receptors at extremely high rate. By con-
12 sidering random walks on graphs, we introduce and analyze a simplified math-
13 ematical model in order to understand this extremely efficient learning process.
14 The structure of the graph reflects the choice of the mutation rule. We focus
15 on the impact of this choice on typical time-scales of the graphs' exploration.
16 We derive explicit formulas to evaluate the expected hitting time to cover a
17 given Hamming distance on the graphs under consideration. This character-
18 izes the efficiency of these processes in driving antibody affinity maturation.
19 In a further step we present a biologically more involved model and discuss
20 its numerical outputs within our mathematical framework. We provide as well
21 limitations and possible extensions of our approach.

22 **Keywords** Random Walks on Graphs · Hypercube · Hitting Time ·
23 Mutational Model · Evolutionary Landscape · Somatic Hypermutation

I. Balelli · V. Milišić
Université Paris 13, Institut Galilée, Département de Mathématiques.
99, avenue Jean-Baptiste Clément 93430 - Villetaneuse - France.
Tel.: +33-1-49-40-35-39 ; +33-1-49-40-35-91
Fax: +33-1-49-40-35-68
E-mail: balelli@math.univ-paris13.fr
E-mail: milisic@math.univ-paris13.fr

G. Wainrib
Ecole Normale Supérieure, Département d'Informatique.
45 rue d'Ulm, 75005 - Paris - France.
E-mail: gilles.wainrib@ens.fr

24 **Contents**

25	1 Introduction	2
26	2 A basic mutational model	5
27	3 More mutational models: how does the structure of the hypercube change?	19
28	4 Modeling issues	33
29	5 Conclusion	46

30 **1 Introduction**

31 Understanding the role and functional implication of mutations is a central
 32 question in biological evolutionary theory [27, 79, 33, 25], but also for the study
 33 of evolutionary algorithms [5, 2]. Beyond the mutation rate, which is natu-
 34 rally an important parameter, our aim in this article is to highlight the role
 35 of various mutation rules on the exploration of the space of traits. In our
 36 mathematical framework, configurations are represented as vertices of a graph
 37 which are connected if there exists a mutation allowing to pass from one trait
 38 to another. We are mainly interested in understanding the characteristic time-
 39 scales for the exploration of the state-space as a function of the mutation rule.
 40 To this end, we relate mutation rules with specific graph topologies and build
 41 upon random walks on graphs and spectral graph theories to analyze resulting
 42 time-scales.

43
 44 More precisely, beyond general theoretical results, we are particularly inter-
 45 ested to apply our framework to the B-cell affinity maturation in Germinal
 46 Centers (GCs). The adaptive immune system is able to create a specific re-
 47 sponse against almost any kind of pathogens penetrating our organism and
 48 inflicting diseases. This task is performed by the production of high affinity
 49 antigen-specific antibodies. These proteins are produced by B-lymphocytes
 50 which are submitted to a learning process improving their affinity to recog-
 51 nize a particular antigen. This process is called Antibody Affinity Maturation
 52 (AAM) and takes place in GCs [54]. Even if substantial progress has been
 53 made in adaptive immunology, since somatic hypermutation was discovered
 54 by the nobel price Susumu Tonegawa [74] in 1987, there are still facts that re-
 55 main unclear about the GC reaction and the exact dynamics of AAM. Indeed,
 56 it seems difficult to make exact measurements of the antigenic repertoire *in*
 57 *vivo* inside a single GC, following and sequencing each B-cell at any time, or
 58 to have precise spatial and temporal data about lymphocytes within the GC
 59 during an immune response, or to understand the exact dynamic of mutation
 60 and selection of B-cells while they are submitted to AAM (*e.g.* [26, 57]). Never-
 61 theless, some refined techniques start to be available [72, 31], showing possible
 62 correlations between proliferation and mutation rates with respect to B-cells'
 63 affinity to the presented antigen. This provides further motivation for setting
 64 appropriate mathematical frameworks to describe such systems.

66 The affinity of a B-cell is biologically observed as a matching between
67 the B-cell receptor (BCR) and the antigen. We aim at understanding how
68 mutation rules allow to explore possible trait-configurations of BCRs. The
69 mutational mechanism that B-cells undergo in GCs to improve their affinity
70 is called Somatic Hypermutation (SHM): it targets, at a very high rate, the
71 DNA encoding for the specific portion of the BCR involved in the binding
72 with the antigen, called Variable (V) region. SHM can introduce mutations
73 at all four nucleotides, and mutation hot-spots have been identified [73,23,
74 71]. The effect of these mutations on the BCR, once expressed on the outer
75 surface of B-cells, is very complex, as the substitution of a single amino-acid
76 can modify the geometrical structure of the BCR, creating or deleting bonds
77 (see [1], Chapter 4, for more details about the crystal structure of BCRs and
78 their binding with antigens).

79
80 Although mutations occur at the level of the DNA, their outcome might be
81 expressed at the level of amino-acids composing the BCR. In the present paper,
82 SHMs are taken in account this way (Section 4.3). However, the structure of
83 our mathematical model can be left substantially unchanged when considering
84 mutations at the DNA level, which leads to modify the definition of affinity
85 and the size of the state-space.

86
87 There already exists a certain number of mathematical models about GC
88 reaction and AAM. In particular, [42,43] proposed deterministic population
89 modeling of SHM and AAM, considering for instance the hypothesis of recycling
90 mechanisms during GC reaction, later investigated by experiments [76].
91 In [56,59,29,36], the authors introduced and discussed several immunological
92 problems, such as the size of the repertoire, or the strength of antigen-antibody
93 binding, or the pourcentage of recycling. They provide suitable mathematical
94 tools, using both deterministic and probabilistic approaches, together with numerical
95 simulations. More recently, biologically very detailed models of GCs
96 were proposed [50,65], using, for instance, agent-based models [51], mostly analyzed
97 through extensive numerical simulations. Our aim here is not to build
98 a very complex model, but rather to contribute to the theoretical foundation
99 of adaptive immunity modeling through the mathematical analysis of generic
100 mutation models on graphs. So far, this approach has not been developed and
101 applied to GC reaction and AAM modeling. In particular, this framework enables
102 the study of various mutation rules, as for instance, affinity-dependent
103 mutations, which are currently debated in the biological literature [31]. Our
104 mathematical framework shares some similarities with the NK models proposed
105 by S. A. Kauffman and E. D. Weinberger in [39], for instance the choice
106 of the hypercube vertex set as the basic structure to define the affinity landscape
107 of BCRs. Nevertheless their approach and goals are fundamentally different
108 from ours. Indeed, in [39] the graph which defines the mutational rule is
109 predefined (*i.e.* they refer only to the basic mutational rule we introduced
110 as well in Section 2), while the affinity function changes according to the main
111 parameters of the model, N and k for instance. Therefore, the random walks

112 over these affinity landscapes, modeling the maturation of the immune re-
113 sponse, are biased with respect to the affinity gradient. In our mathematical
114 framework the structure of the graph reflects the mutational rule, hence it is
115 not predefined. Moreover, since in this paper we only take into account mu-
116 tations, the random walks over the state-space are not biased by the fitness
117 of each trait to the target one. From our point of view the selection pressure
118 should be taken into account as a separate operator (see below).

119
120 This research is also motivated by important biotechnological applica-
121 tions. The fundamental understanding of the evolutionary mechanisms in-
122 volved in antibody affinity maturation have been inspiring many methods for
123 the synthetic production of specific antibodies for drugs, vaccines or cancer im-
124 munotherapy [4, 45, 67]. Indeed, this production process involves the selection
125 of high affinity peptides and requires smart methods to generate an appropriate
126 diversity [18]. Beyond the biomedical motivations, the study of this learning
127 process has also given rise in recent years to a new class of bio-inspired algo-
128 rithms such as in [16, 58], mainly addressed to solve optimization and learning
129 problems [13].

130
131 In this article, we consider pure mutational models obtained as random
132 walks on graphs given by alterations of the edge set of the N -dimensional
133 hypercube. We focus on the variation of hitting times as a function of the un-
134 derlying graphs, hence relating mutation rules to the characteristic time-scales
135 of the process. Our intention here is not to provide biologically relevant out-
136 comes, since the AAM involves several mechanisms (division, selection, etc)
137 that we do not take into account in this article. Instead we provide a rigorous
138 analysis of an essential single building block: mutation. We study the structure
139 of RWs on the hypercube and compute hitting times depending on the graph
140 associated to the mutational rule. We prove that they are proportional to the
141 number of vertices (see Table 2). Therefore our specific approach consists in
142 observing how different mutational rules allow to explore the state-space and
143 lead a naive B-cell to build the fittest possible trait. We are not interested
144 here in proposing new statistical or phylogenetic strategies to infer the more
145 realistic phylogenetic trees given a final antibodies repertoire [30, 17]. Nev-
146 ertheless we define accurately the biological context since it is relevant for
147 further steps. Clearly, other mechanisms such division and mutations provide
148 significant biases of hitting times, our approach consists in studying precisely
149 the differences when enriching our model with supplementary bricks. For in-
150 stance, by branching we introduce a population dispatched on the vertices of
151 the hypercube which decreases the hitting time, but at the cost of the bio-
152 logical maintaining of the population [6]. This is our strategy here and in the
153 forthcoming papers [6, 7].

154
155 Section 2 contains results on random walks theory [55, 52, 61] and, more
156 specifically, random walks on graphs [49, 3]. This is a topic of active research
157 due to the great number of important applications in recent years, such as

158 graph clustering [64], ranking algorithms for search-engines [10,37], or social
 159 network modeling [41,32,44]. We start with the most basic mutational model
 160 which is the simple random walk on the N -dimensional hypercube [22,34,21,
 161 77]. We set notations in order to define the models, then we overview various
 162 properties of random walks on graphs, and establish particular results in the
 163 case of the hypercube. In Section 3 we study several mutation rules and their
 164 effects on the structure of the graph and, consequently, its associated random
 165 walk. In particular we compute the hitting times: starting from a random initial
 166 condition, we count the expected time to reach the target node with the best
 167 fitness. We use both spectral and probabilistic methods. We especially focus on
 168 two mutation rules that are the combination of simpler ones: the class switch
 169 of 1 or 2-length strings (Section 3.1.3), where the mutation rule depends on the
 170 distance to the target, and the mutation rule which allows to do more than a
 171 single mutation at each step (Section 3.1.4). Table 2 in Section 3.2 summarizes
 172 the main results of Section 2 and 3: we display expected times to reach some
 173 position of the graph, as a function of each mutation rule. Finally, Section 4 is
 174 dedicated to modeling aspects and discussions about possible extensions and
 175 limitations of the proposed framework.

176 2 A basic mutational model

177 In this section we set the general mathematical framework, which we keep in
 178 order to pattern and study mutational mechanisms discussed in the current
 179 section and in Section 3. Indeed, we state a basic mutational model. The choice
 180 of this environnement is motivated by the modeling of amino-acids chains and
 181 their modifications during SHM. It is for this reason that we often recall bio-
 182 logical facts and refer to BCRs and antigens. Nevertheless, this framework is
 183 flexible and adapts to different mutational rules in a more general evolutionary
 184 context.

185
 186 We assume that it is possible to classify the amino-acids into 2 classes
 187 denoted by 0 and 1 respectively (they could represent amino-acids negatively
 188 and positively charged respectively). Henceforth BCRs and antigen are repre-
 189 sented by binary strings of same fixed length N , hence, the state-space of all
 190 possible BCR configurations is $\{0,1\}^N$. We will give some more details about
 191 these hypotheses in Section 4.3.

192 **Definition 1** We denote by \mathcal{H}_N the standard N -dimensional hypercube. BCR
 193 and antigen configurations are represented by vertices of \mathcal{H}_N , denoted by \mathbf{x}_i
 194 with $1 \leq i \leq 2^N$, or sometimes simply by their indices. We denote the antigen
 195 target vertex by $\bar{\mathbf{x}}$: it is given at the beginning of the process and never changes.

196 We suppose that there is a single B-cell entering the GC reaction. The
 197 configuration of its receptors is denoted by \mathbf{X}_0 . If \mathbf{X}_t is the configuration of
 198 the BCR after t mutations, then depending on the mutational rule, one or more
 199 bits in \mathbf{X}_t can change after the next mutation. This gives rise to a Random

200 Walk (RW) on $\{0, 1\}^N$, where a mutation on the BCR corresponds to a jump
 201 to a neighbor node. Of course, the definition of neighbors changes depending
 202 on the mutation rules we introduce (we specify the neighborhood set each time
 203 we discuss a new mutation rule). In a general way:

204 **Definition 2** Given $\mathbf{x}_i, \mathbf{x}_j \in \{0, 1\}^N$, we say that \mathbf{x}_i and \mathbf{x}_j are neighbors,
 205 and denote $\mathbf{x}_i \sim \mathbf{x}_j$, if there exists at least one edge (or loop) between them.

206 As far as the complementarity is concerned, we have to make a further
 207 simplification. As we have already discussed in the Introduction, the tridimen-
 208 sional structure of the BCR is hard to model. For this reason we consider a
 209 linear contact, *i.e.* positively charged amino-acids are complementary to neg-
 210 atively charged ones when they are at the same position within the binary
 211 string. For the sake of simplicity, we state that 0 matches with 0 and 1 with
 212 1 (we can suppose that the antigen representing string is given in its comple-
 213 mentary form). Formally, we define the affinity as the number of identical bits
 214 shared by the BCR representing string and $\bar{\mathbf{x}}$. Equivalently, one can see $\bar{\mathbf{x}}$ as
 215 the optimal BCR trait, with the highest affinity for the immunizing antigen.

216 **Definition 3** For all $\mathbf{x}_i \in \{0, 1\}^N$, its affinity with $\bar{\mathbf{x}}$, $\text{aff}(\mathbf{x}_i, \bar{\mathbf{x}})$ is given by
 217 $\text{aff}(\mathbf{x}_i, \bar{\mathbf{x}}) := N - h(\mathbf{x}_i, \bar{\mathbf{x}})$, where $h(\cdot, \cdot) : \{0, 1\}^N \times \{0, 1\}^N \rightarrow \{0, \dots, N\}$ returns
 218 the Hamming distance.

Definition 4 For all $\mathbf{x} = (x_1, \dots, x_N), \mathbf{y} = (y_1, \dots, y_N) \in \{0, 1\}^N$, their Ham-
 ming distance is given by:

$$h(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^N \delta_i \quad \text{where} \quad \delta_i = \begin{cases} 1 & \text{if } x_i \neq y_i \\ 0 & \text{otherwise} \end{cases}$$

219 Other definitions of affinity are often (*e.g.* [50]) constructed as functions
 220 of the Hamming distance $\text{aff}(\mathbf{x}_i, \bar{\mathbf{x}}) = F(h(\mathbf{x}_i, \bar{\mathbf{x}}))$, for instance with F given
 221 by the Gaussian probability density function. These modeling aspects become
 222 important when considering the selection mechanism, which is not treated in
 223 the present article. Therefore, for our purpose, we can focus on the above def-
 224 inition of affinity.

225 As a first basic mutational rule, we study single switch-type mutations: at
 226 each time step a randomly chosen amino-acid within the BCR binary string
 227 switches its amino-acid class. This clearly leads us to a Simple Random Walk
 228 (SRW) on \mathcal{H}_N . Indeed, we formalize it as follows:

230 **Definition 5** Let $\mathbf{X}_n \in \mathcal{H}_N$ be the BCR at step n . Let $i \in \{1, \dots, N\}$ be a ran-
 231 domly chosen index. Then $\mathbf{X}_{n+1} := (X_{n,1}, \dots, X_{n,i-1}, 1 - X_{n,i}, X_{n,i+1}, \dots, X_{n,N})$.

232 *Remark 1* Referring to Definition 2 of neighborhood, as we consider here the
 233 standard N -dimensional hypercube, $\forall \mathbf{x}_i, \mathbf{x}_j \in \mathcal{H}_N, \mathbf{x}_i \sim \mathbf{x}_j \Leftrightarrow h(\mathbf{x}_i, \mathbf{x}_j) = 1$.

We denote the transition probability matrix of the SRW on \mathcal{H}_N by \mathcal{P}_N or simply by \mathcal{P} if no misunderstanding is possible. For all $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{H}_N$:

$$\mathbb{P}(\mathbf{X}_n = \mathbf{x}_j | \mathbf{X}_{n-1} = \mathbf{x}_i) =: p(\mathbf{x}_i, \mathbf{x}_j) = \begin{cases} 1/N & \text{if } \mathbf{x}_j \sim \mathbf{x}_i, \\ 0 & \text{otherwise.} \end{cases}$$

234 The entries of \mathcal{P} are $(p(\mathbf{x}_i, \mathbf{x}_j))_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{H}_N}$. The unique stationary distribution
 235 for \mathcal{P} is the homogeneous probability distribution on \mathcal{H}_N , denoted by $\boldsymbol{\pi}$:
 236 $\forall \mathbf{x}_i \in \mathcal{H}_N, \pi_i := \boldsymbol{\pi}(\mathbf{x}_i) = 2^{-N}$. Indeed, $(\mathbf{X}_n)_{n \geq 0}$ is clearly reversible with respect to $\boldsymbol{\pi}$. The uniqueness follows by the Ergodic Theorem.

238
 239 We also recall a property of \mathcal{H}_N that we will have to deal with: the bipar-
 240 titeness.

241 **Definition 6** A graph $G = (V, E)$ is bipartite if there exists a partition of the
 242 vertex set $V = V_1 \sqcup V_2$, s.t. every edge connects a vertex in V_1 to a vertex in
 243 V_2 .

244 Typically a bipartition of the hypercube can be obtained by separating the
 245 vertices with an odd number of 1's in their string from those with an even
 246 number of 1's. In Figure 1 we emphasize the bipartite structure of the hyper-
 247 cube \mathcal{H}_3 .

248

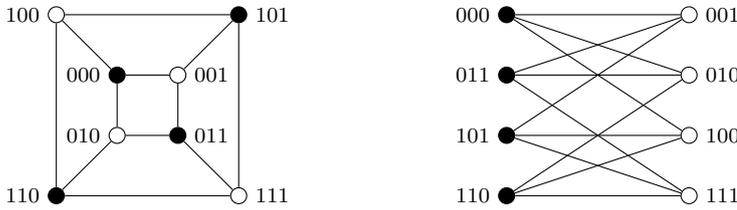


Figure 1: Hypercube for $N = 3$ showing its bipartite structure.

249 A direct and elementary consequence of this property is the periodic be-
 250 havior of the SRW on \mathcal{H}_N , which in particular causes some problems for the
 251 convergence through $\boldsymbol{\pi}$. This problem is classically overcome by adding N
 252 loops at each vertex, that makes this RW become a *lazy Markov chain* [48].
 253 The corresponding transition probability matrix is given by $\mathcal{P}_L := (\mathcal{P} + I_{2^N})/2$,
 254 where I_n denotes the n -dimensional identity matrix.

255 2.1 Spectral analysis

256 Most matrices describing the characteristics of the SRW on \mathcal{H}_N can be ob-
 257 tained recursively, thanks to the recursive construction of the hypercube and
 258 the operation of cartesian product between two graphs.

259 **Definition 7** Given two graphs $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$, the cartesian
 260 product between G_1 and G_2 , $G_1 \times G_2$, is a graph with vertex set $V = V_1 \times V_2 =$
 261 $\{(u, v) \mid u \in V_1, v \in V_2\}$. Two different vertices (u_1, v_1) and (u_2, v_2) are adjacent
 262 in $G_1 \times G_2$ if either $u_1 = u_2$ and $v_1 v_2 \in E_2$ or $v_1 = v_2$ and $u_1 u_2 \in E_1$.

263 It is a known result [34] that for $N > 1$, \mathcal{H}_N is obtained from \mathcal{H}_{N-1} as:
 264 $\mathcal{H}_N = \mathcal{H}_{N-1} \times \mathcal{H}_1$. This characteristic implies the recursive construction of the
 265 adjacency matrix and allows to determine the corresponding eigenvalues and
 266 eigenvectors. We denote by A_N the adjacency matrix corresponding to \mathcal{H}_N ;
 267 by I_n the n -dimensional identity matrix. Then we have:

$$A_1 = \begin{array}{c} 0 \\ 1 \end{array} \begin{array}{c|c} 0 & 1 \\ \hline 1 & 0 \end{array}; \quad A_2 = \begin{array}{cc} 00 & \begin{array}{c|c} 0 & 1 \\ \hline 1 & 0 \end{array} \\ 01 & \begin{array}{c|c} 1 & 0 \\ \hline 0 & 1 \end{array} \\ 10 & \begin{array}{c|c} 1 & 0 \\ \hline 0 & 1 \end{array} \\ 11 & \begin{array}{c|c} 0 & 1 \\ \hline 1 & 0 \end{array} \end{array} = \begin{array}{c|c} A_1 & I_2 \\ \hline I_2 & A_1 \end{array}$$

268 Here we wrote in gray the strings corresponding to each row: in order to obtain
 269 the adjacency matrices in this form, we simply have to order vertices of \mathcal{H}_N
 270 in lexicographical order.

271

By iteration we obtain [28]:

$$A_n = \begin{array}{c|c} A_{n-1} & I_{2^{n-1}} \\ \hline I_{2^{n-1}} & A_{n-1} \end{array}$$

272 This iterative construction allows also to determine recursively the spectra
 273 of A_N and, consequently, of $\mathcal{P}_N = A_N/N$ (as \mathcal{H}_N is a N -regular graph, the
 274 transition probability matrix corresponds to the adjacency matrix divided by
 275 N). Here below we recall the explicit values of the eigenvalues of A_N and \mathcal{P}_N
 276 respectively. An extensive proof can be found in [28].

277 **Theorem 1** *The eigenvalues of A_N are: $N, N-2, N-4, \dots, -N+4, -N+2, -N$. If we order the $N+1$ distinct eigenvalues of A_N as $\lambda_1^A > \lambda_2^A > \dots >$
 278 λ_{N+1}^A , then the multiplicity of λ_k^A is $\binom{N}{k-1}$, $1 \leq k \leq N+1$*

280 **Corollary 1** *The eigenvalues of \mathcal{P}_N are: $1, 1-2/N, 1-4/N, \dots, -1+4/N, -1+2/N, -1$. If we order the $N+1$ distinct eigenvalues of \mathcal{P} as $\lambda_1 > \lambda_2 > \dots >$
 281 λ_{N+1} , then the multiplicity of λ_k is $\binom{N}{k-1}$, $1 \leq k \leq N+1$*

Finally we recall the expression of the eigenvectors of A_N (and then also of \mathcal{P}), that we gather together into a matrix. The eigenvectors for A_1 are:

$$\mathbf{z}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \text{ for } \lambda_1^A = 1 \quad \text{and} \quad \mathbf{z}_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \text{ for } \lambda_2^A = -1 \Rightarrow \mathcal{Z}_1 = [\mathbf{z}_1, \mathbf{z}_2]$$

Thanks to the relations between the cartesian product of two graphs and their eigenvectors, it follows by induction that [28]:

$$\mathcal{Z}_n = \left(\begin{array}{c|c} \mathcal{Z}_{n-1} & \mathcal{Z}_{n-1} \\ \hline \mathcal{Z}_{n-1} & -\mathcal{Z}_{n-1} \end{array} \right)$$

283 Finally, one renormalizes each vector \mathbf{z}_i multiplying it by $\sqrt{2^{-N}}$. We denote
 284 by Q_N the resulting matrix, where each column is a 2^N vector $\mathbf{v}_i = \sqrt{2^{-N}}\mathbf{z}_i$.

285 2.2 Evolution of Hamming distances to a fixed node

286 In this section we focus on the distance process, which is the process obtained
 287 from the SRW on \mathcal{H}_N by looking at the Hamming distance between the B-cell
 288 representing string at each mutation step and the antigen target representing
 289 string. More precisely, $(D_n)_{n \geq 0} := (h(\mathbf{X}_n, \bar{\mathbf{x}}))_{n \geq 0}$ is a RW on $\{0, \dots, N\}$. From
 290 a biological point of view this process represents the evolution of the affinity
 291 of the mutating B-cell to the presented antigen. The idea of analyzing the distance
 292 of a RW on a graph to some position, where distance means the minimal
 293 number of steps that separate two positions, is not unusual. N. Berestycki in
 294 [9] applied that to genome rearrangements, where the distance on the graph
 295 corresponds biologically to the minimal number of reversals or other mutations
 296 needed to transform one genome into the other. Due to the perfect symme-
 297 try of the graph under consideration and our particular choice of the affinity
 298 (which is directly related to the Hamming distance), by studying (D_n) we
 299 reduce considerably the number of vertices, passing from 2^N to $N + 1$ nodes,
 300 without losing the most important properties of the corresponding transition
 301 matrix. However, if we consider more complicated models of mutation, it is
 302 not possible to reduce the study of the process to the distances to a fixed node.
 303 In Figure 2 we show explicitly how to pass from (\mathbf{X}_n) to (D_n) : since $\bar{\mathbf{x}}$ is fixed
 304 and known, we are able to group the vertices by their Hamming distance to $\bar{\mathbf{x}}$.
 305 Moreover we keep the original probability of going to the next distance class
 306 by considering weighted and directed edges.

307
 308 The transition probability matrix for (D_n) , denoted by \mathcal{Q} , is given by
 309 Proposition 1 below.

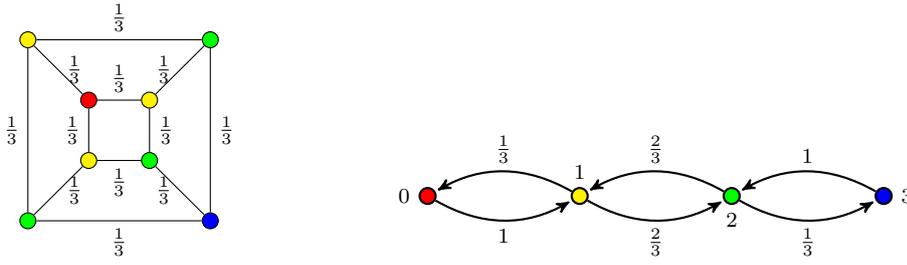


Figure 2: From the (\mathbf{X}_n) process (on the left) to the (D_n) process (on the right) (case $N = 3$). Near each arrow the probability to travel in the corresponding direction is exhibited. The red vertex always corresponds to \bar{x} , while we represent vertices at the same distance with the same color (yellow for $h = 1$, green for $h = 2$, and blue for $h = 3$).

310 **Proposition 1** For all $d, d' \in \{0, \dots, N\}$:

$$\mathbb{P}(D_n = d' \mid D_{n-1} = d) =: q(d, d') = \begin{cases} d/N & \text{if } d' = d-1 \\ (N-d)/N & \text{if } d' = d+1 \\ 0 & \text{if } |d' - d| \neq 1 \end{cases} \quad (1)$$

311 $\mathcal{Q} = (q(d, d'))_{d, d' \in \{0, \dots, N\}}$ is a $(N+1) \times (N+1)$ tridiagonal matrix where
 312 the main diagonal consists of zeros. The stationary distribution for \mathcal{Q} is the
 313 binomial probability distribution $\mathcal{B}(N, \frac{1}{2}) = \left(C_N^d \frac{1}{2^N}\right)_{d \in \{0, \dots, N\}}$, where $C_N^d =$
 314 $\binom{N}{d} = \frac{N!}{d!(N-d)!}$ is the binomial coefficient. It is the unique stationary distribu-
 315 tion for \mathcal{Q} : a simple calculation points out the fact that $(D_n)_{n \geq 0}$ is reversible
 316 with respect to $\mathcal{B}(N, \frac{1}{2})$, then the uniqueness follows by the Ergodic Theorem.

317
 318 Anew, we have to deal with bipartiteness: the graph we are taking into
 319 account in this section is clearly bipartite, since we can separate its vertices
 320 into two subsets containing odd and even nodes respectively and no edge
 321 connects any vertices in the same subset. In order to overcome this problem
 322 we add N loops at each vertex $\mathbf{x}_i \in \mathcal{H}_N$ which means that the new transition
 323 probability matrix for the (D_n) process is, for all $d, d' \in \{0, \dots, N\}$:

$$\mathbb{P}(D_n = d' \mid D_{n-1} = d) =: q_L(d, d') = \begin{cases} 1/2 & \text{if } d' = d \\ d/(2N) & \text{if } d' = d-1 \\ (N-d)/(2N) & \text{if } d' = d+1 \\ 0 & \text{if } |d' - d| \neq 1 \end{cases} \quad (2)$$

324 We denote by $\mathcal{Q}_L := (q_L(d, d'))_{d, d' \in \{0, \dots, N\}}$.

Proposition 2 $(D_n)_{n \geq 0}$ converges in law to a binomial random variable with parameters N and $1/2$. Explicitly:

$$(\mathcal{Q}_L)_d \rightarrow \mathcal{B}\left(N, \frac{1}{2}\right)_d \quad \text{for } n \rightarrow +\infty$$

Proof The proof follows directly observing that \mathcal{Q}_L represents an irreducible and, now, aperiodic MC, with the same stationary distribution as \mathcal{Q} (see [55] for a proof of the general result). \square

325 The spectral analysis of \mathcal{Q} gives the following result.

326 **Theorem 2** For fixed N , the spectra of the transition probability matrix \mathcal{Q}
 327 corresponding to the (D_n) process is composed by the same $N + 1$ distinct
 328 eigenvalues as the spectra of \mathcal{P} , each with multiplicity 1.

Proof The proof consists of a simple calculation of the eigenvalues of matrix \mathcal{Q} , which is easily done for $N = 1, 2$. Then we reason by iteration. We can also give the system we use for determining the eigenvectors. For fixed N let us denote by $\lambda_{\pm k}$ the eigenvalue $\frac{\pm(N-2k)}{N}$ for $0 \leq k \leq \lfloor N/2 \rfloor$. We denote by $\mathbf{x}_{\pm k}$ the corresponding unknown eigenvector. Then we have the following matrix equation:

$$\mathcal{Q}\mathbf{x}_{\pm k} = \lambda_{\pm k}\mathbf{x}_{\pm k}$$

Which is:

$$\left\{ \begin{array}{l} x_{\pm k,2} = \lambda_{\pm k}x_{\pm k,1} \\ \frac{1}{N}x_{\pm k,1} + \frac{N-1}{N}x_{\pm k,3} = \lambda_{\pm k}x_{\pm k,2} \\ \frac{2}{N}x_{\pm k,2} + \frac{N-2}{N}x_{\pm k,4} = \lambda_{\pm k}x_{\pm k,3} \\ \vdots \\ \frac{N-1}{N}x_{\pm k,N-1} + \frac{1}{N}x_{\pm k,N+1} = \lambda_{\pm k}x_{\pm k,N} \\ x_{\pm k,N} = \lambda_{\pm k}x_{\pm k,N+1} \end{array} \right.$$

\square

Remark 2 Using the classical results of S. N. Ethier and T. G. Kurtz [24] it is possible to prove that, denoting by $x_N(t)$ the process $x_N(t) = \frac{D_{\lfloor Nt \rfloor}}{N}$,

it converges in probability through $x(t)$, solution of the differential equation $\dot{x}(t) = -2x(t) + 1$ on a finite time window:

$$\forall \varepsilon > 0, \forall T > 0, \mathbb{P} \left(\sup_{t \in [0, T]} |x_N(t) - x(t)| > \varepsilon \right) \rightarrow 0 \quad \text{for } N \rightarrow \infty.$$

329 *Remark 3* We can easily observe that $x(t)$ rapidly converges to $1/2$ for all
 330 $x_0 \in [0, 1]$. In particular if we start at $x_0 = 1/2$, we stay there for all t . That
 331 suggests that the (D_n) process, for N going to infinity, reaches a value of about
 332 $N/2$ exponentially fast, and then tends to remain there.

From an heuristic viewpoint we can explain how we derived the above equation. First of all, we take into account the following rescaled process:

$$x_n := D_n/N$$

As $(D_n) \in \{0, \dots, N\}$, $x_n \in [0, 1]$. Denoting by $q_n(x) = \mathbb{P}(x_n = x)$ and using Equation (1), we have:

$$q_{n+1}(x) = (1-x)q_n \left(x - \frac{1}{N} \right) + xq_n \left(x + \frac{1}{N} \right)$$

Now we apply the Taylor theorem for $N \gg 1$:

$$q_{n+1}(x) = (1-x) \left(q_n(x) - \frac{1}{N} q'_n(x) + o \left(\frac{1}{N} \right) \right) + x \left(q_n(x) + \frac{1}{N} q'_n(x) + o \left(\frac{1}{N} \right) \right)$$

From which we get:

$$q_{n+1}(x) - q_n(x) = \frac{1}{N} (x - (1-x)) q'_n(x) + o \left(\frac{1}{N} \right)$$

Defining the process $\tilde{q}(t, x) = q_{\lfloor Nt \rfloor}(x)$, with $t = \frac{n}{N}$, we obtain:

$$\partial_t \tilde{q}(t, x) = (2x - 1) \partial_x \tilde{q}(t, x) + o \left(\frac{1}{N} \right)$$

333 And consequently, the corresponding transport equation is:

$$\partial_t q(t, x) = (2x - 1) \partial_x q(t, x) \tag{3}$$

The differential equation associated with Equation (3) (its characteristic equation) is:

$$\dot{x}(t) = -2x(t) + 1$$

which has solution:

$$x(t) = \frac{1}{2} + \left(x_0 - \frac{1}{2} \right) e^{-2t}$$

334 It is also possible to derive a diffusion approximation by expanding the gener-
 335 ator at second order.

336 2.3 Hitting times

337 In this section we give explicit formulas to compute the hitting time from node
 338 \mathbf{x}_i to \mathbf{x}_j : the expected number of steps before \mathbf{x}_j is visited, starting from \mathbf{x}_i .
 339 More precisely, we define by $\tau_{\{\mathbf{x}_j\}} := \inf\{n \geq 0 \mid \mathbf{X}_n = \mathbf{x}_j\}$: we are interested
 340 in studying its expectation, $\mathbb{E}_{\mathbf{x}_i}[\tau_{\{\mathbf{x}_j\}}]$. The formula we gave in Section 2.3.1
 341 is directly obtained from the more general one given by L. Lovász in [49]: we
 342 recall it simply because we will need it later. On the other hand, the formula
 343 given in Section 2.3.2 is obtained from the (D_n) process and the procedure is
 344 inspired by those used in [47].

345 2.3.1 Analysis of $\mathbb{E}_{\mathbf{x}_0}[\tau_{\{\bar{\mathbf{x}}\}}]$ using the spectrum of \mathcal{P} .

346 **Definition 8** Let H be the $2^N \times 2^N$ symmetric matrix having as $(i, j)^{\text{th}}$ entry:
 347 $(H)_{ij} = H(i, j) = \mathbb{E}_{\mathbf{x}_i}[\tau_{\{\mathbf{x}_j\}}]$ for all $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{H}_N$. Clearly $H(i, i) = 0$ for all i .

348 The N -regularity of the graph implies that:

$$H(i, j) = 1 + \sum_{\{k \mid h(i, k)=1\}} \mathcal{P}_{ik} H(k, j) = 1 + \frac{1}{N} \sum_{\{k \mid h(i, k)=1\}} H(k, j) \quad \text{for } i \neq j \quad (4)$$

To relate the hitting time with the spectrum, we first define $F := J + \mathcal{P}H - H$, where J is a $2^N \times 2^N$ matrix whose entries are all 1. From Equation (4), it follows that F is a diagonal matrix, as $(H)_{ij} = (J)_{ij} + (\mathcal{P}H)_{ij}$ for $i \neq j$. Moreover $F'\pi = \mathbf{1}$, where $\mathbf{1} = (1, \dots, 1)'$, since

$$F'\pi = (J + (\mathcal{P} - I_{2^N})H)'\pi = J\pi + H'(\mathcal{P} - I_{2^N})'\pi = J\pi + H'(\mathcal{P}'\pi - \pi) = J\pi = \mathbf{1}$$

349 Therefore, we deduce that $F = 2^N I_{2^N}$ and H is solution of

$$(I_{2^N} - \mathcal{P})H = J - 2^N I_{2^N} \quad (5)$$

350

351 **Theorem 3** Given a SRW on \mathcal{H}_N , the hitting time from vertex i to j is given
 352 by:

$$H(i, j) = 2^N \sum_{k=2}^{2^N} \frac{1}{1 - \lambda_k} (v_{kj}^2 - v_{ki}v_{kj}), \quad (6)$$

353 where λ_k is the k^{th} eigenvalue of \mathcal{P} and v_{ki} corresponds to the i^{th} component
 354 of the k^{th} eigenvector of \mathcal{P} , as given in Section 2.1.

Proof We can not directly solve equation (5), since matrix $(I_{2^N} - \mathcal{P})$ is singular. The spectral decomposition theorem insures that $\mathbb{R}^{2^N} = \oplus_{i=1}^{2^N} \text{Span}\{\mathbf{v}_i\}$. On the subspace $\oplus_{i=2}^{2^N} \text{Span}\{\mathbf{v}_i\}$, $(I_{2^N} - \mathcal{P})$ is invertible. At the same time, the

right hand side in (5) reduces to a constant times the identity matrix when restricted to this same subspace. Thus a possible candidate solving (5) is:

$$\tilde{H} = -2^N \sum_{i=2}^{2^N} (1 - \lambda_i)^{-1} \mathbf{v}_i \mathbf{v}_i'$$

Nevertheless, for every vector $\mathbf{w} \in \mathbb{R}^{2^N}$, $\tilde{H} + \mathbf{1}\mathbf{w}'$ is a solution of (5) as well. Thus H can be unambiguously determined by imposing the condition over its main diagonal: $H(i, i) = 0$ for all $i \in \{0, \dots, 2^N\}$. \square

355 *2.3.2 Analysis of $\mathbb{E}_{\mathbf{x}_0}[\tau_{\{\bar{\mathbf{x}}\}}]$ from the D_n viewpoint.*

356 For the sake of simplicity, we denote $H(D_0) := \mathbb{E}_{\mathbf{x}_0}[\tau_{\{\bar{\mathbf{x}}\}}]$ as it depends only
357 on the initial Hamming distance of \mathbf{X}_0 to $\bar{\mathbf{x}}$, D_0 .

Remark 4 Due to (1), starting at point \mathbf{x}_0 with $D_0 = \bar{d}$, we have:

$$\begin{cases} \mathbb{P}(D_1 = \bar{d} + 1 | D_0 = \bar{d}) =: q(\bar{d}, \bar{d} + 1) = (N - \bar{d})/N \\ \mathbb{P}(D_1 = \bar{d} - 1 | D_0 = \bar{d}) =: q(\bar{d}, \bar{d} - 1) = \bar{d}/N \end{cases}$$

358 We are now able to define a new recursive formula for (4), which will be more
359 convenient if evaluated explicitly:

$$H(\bar{d}) = 1 + \frac{N - \bar{d}}{N} H(\bar{d} + 1) + \frac{\bar{d}}{N} H(\bar{d} - 1) \quad (7)$$

360 with boundary conditions:

$$H(0) = 0 \text{ and } H(1) = 2^N - 1 = \sum_{j=0}^N C_N^j - 1 \quad (8)$$

Taking the difference $\Delta(\bar{d}) := H(\bar{d}) - H(\bar{d} - 1)$, we obtain:

$$\Delta(\bar{d} + 1) = H(\bar{d} + 1) - H(\bar{d}) = \frac{\bar{d}}{N} (\Delta(\bar{d} + 1) + \Delta(\bar{d})) - 1$$

361 And finally:

$$\Delta(\bar{d} + 1) = \frac{\bar{d}}{N - \bar{d}} \Delta(\bar{d}) - \frac{N}{N - \bar{d}} \quad \text{with } \Delta(1) = H(1) \quad (9)$$

362 Then we can prove by iteration the following result:

363 **Theorem 4** *Given a SRW on \mathcal{H}_N , the hitting time to cover a Hamming dis-*
364 *tance equal to \bar{d} , $H(\bar{d})$ with $0 \leq \bar{d} \leq N$ is obtained as:*

$$H(\bar{d}) = \sum_{d=0}^{\bar{d}-1} \frac{\sum_{j=1}^{N-1-d} C_N^{d+j} + 1}{C_{N-1}^d} \quad (10)$$

365 *Proof* One have to prove that:

$$\Delta(d+1) = \frac{\sum_{j=1}^{N-1-d} C_N^{d+j} + 1}{C_{N-1}^d} \tag{11}$$

366

$$\begin{aligned} \Delta(d+1) &= \frac{d \cdot \Delta(d)}{N-d} - \frac{N}{N-d} = \frac{d}{N-d} \left(\frac{(d-1) \cdot \Delta(d-1)}{N-(d-1)} - \frac{N}{N-(d-1)} \right) - \frac{N}{N-d} \\ &= \frac{d(d-1) \cdot \Delta(d-1)}{(N-d)(N-(d-1))} - N \left(\frac{d}{(N-d)(N-(d-1))} + \frac{1}{N-d} \right) \end{aligned} \tag{12}$$

367 Proceeding by iteration we obtain two terms, where the first one multiplies
 368 $\Delta(1)$. From Equation (9) we know that $\Delta(1) = H(1) = \sum_{j=0}^N C_N^j - 1$. A convenient
 369 use of the properties of the factorial operator allows us to reach the
 370 following expression:

$$\begin{aligned} (12) &= \frac{d!(N-1-d)!}{(N-1)!} \left(\sum_{j=0}^N C_N^j - 1 \right) - N \left(\frac{d!(N-1-d)!}{(N-1)!} + \frac{d!(N-1-d)!}{2!(N-2)!} + \dots \right. \\ &\quad \left. + \frac{d!(N-1-d)!}{(d-1)!(N-(d-1))!} + \frac{d!(N-1-d)!}{d!(N-d)!} \right) = \\ &= \frac{d!(N-1-d)!}{(N-1)!} \left(1 + \sum_{j=1}^{N-1-d} \frac{N!}{(d+j)!(N-(d+j))!} \right) = \frac{\sum_{j=1}^{N-1-d} C_N^{d+j} + 1}{C_{N-1}^d} \end{aligned}$$

By using again (9), we can now easily express $H(\bar{d})$ in the following way

$$H(\bar{d}) = \sum_{d=0}^{\bar{d}-1} \Delta(d+1) = \sum_{d=0}^{\bar{d}-1} \frac{\sum_{j=1}^{N-1-d} C_N^{d+j} + 1}{C_{N-1}^d}$$

which can be evaluated for reasonable values of N . □

We can immediately observe that $H(\bar{d})$ is a monotonically increasing function. Moreover, H is concave. Indeed, thanks to Proposition 4 we can prove that $\forall d \in \{1, \dots, N-1\}$:

$$H(d) - H(d-1) \geq H(d+1) - H(d) \iff \Delta(d) \geq \Delta(d+1)$$

371 Furthermore, we can evaluate the following limit:

$$\lim_{N \rightarrow \infty} \frac{H(\alpha N)}{2^N} \quad \text{for } \alpha \in]0, 1]. \tag{13}$$

372

373 *Remark 5* The case $\alpha = 0$ is trivial: if $\alpha = 0$ this limit is equal to 0 since
 374 $H(0) = 0$.

375 *Remark 6* Proposition 3 below, which evaluates (13), confirms the statement
 376 made in Remark 3: as N goes to infinity, (D_n) goes quickly to $N/2$ and then
 377 $H(d)$ is always of order $\sim 2^N$ irrespective of $d \neq 0$.

Proposition 3 For all $\alpha \in]0, 1[$:

$$\lim_{N \rightarrow \infty} \frac{H(\alpha N)}{2^N} = 1$$

Proof Since H is an increasing function and by using Equation (10) we have:

$$2^N - 1 = H(1) \leq H(\alpha N) \leq H(N) = \sum_{d=0}^{N-1} \frac{1}{C_{N-1}^d} + \sum_{d=0}^{N-1} \sum_{j=1}^{N-1-d} \frac{C_N^{d+j}}{C_{N-1}^d} =: S_1 + S_2$$

378 We examine the two terms of the last member separately.

$$S_1 \leq 2 + \frac{2}{N-1} + (N-4) \frac{2}{(N-1)(N-2)} \quad (14)$$

379 We can prove it just by looking at Pascal's triangle.

380

Now, if we consider S_2 , we see that there is no contribution for $d = N - 1$, as the internal sum is zero valued. Moreover we have:

$$\sum_{j=1}^{N-1-d} C_N^{d+j} \leq \sum_{j=0}^N C_N^j = 2^N$$

And so:

$$S_2 \leq 2^N \sum_{d=0}^{N-2} \frac{1}{C_{N-1}^d} \stackrel{(14)}{\leq} 2^N \left(1 + \frac{2}{N-1} + (N-4) \frac{2}{(N-1)(N-2)} \right)$$

By putting together all these inequalities and dividing by factor 2^N we get that:

$$1 - \frac{1}{2^N} \leq \frac{H(\alpha N)}{2^N} \leq 1 + \frac{2}{N-1} + \frac{2(N-4)}{(N-1)(N-2)} + \frac{1}{2^N} \left(2 + \frac{2}{N-1} + \frac{2(N-4)}{(N-1)(N-2)} \right)$$

The result comes directly by applying the squeeze theorem. \square

381 This result can be extended to a SRW on a generic state-space \mathcal{S}^N , with
 382 $|\mathcal{S}| = s$. More precisely, one can prove in a similar way as we did for \mathcal{H}_N the
 383 following result:

384 **Proposition 4** The order of magnitude of the hitting time for a switch-type
 385 mutational model on the state-space \mathcal{S}^N , with $|\mathcal{S}| = s$, is s^N , for N big enough.

386 This is the consequence of Theorem 5 and Proposition 5 below.

387 **Theorem 5** Given a SRW on S^N , the hitting time to cover a Hamming distance equal to \bar{d} , $H^s(\bar{d})$ with $0 \leq \bar{d} \leq N$ is obtained as:

$$H^s(\bar{d}) = \sum_{d=0}^{\bar{d}-1} \frac{\sum_{j=d+1}^N C_N^j (s-1)^j}{C_{N-1}^d (s-1)^d} \quad (15)$$

388 **Proposition 5** For all $\alpha \in]0, 1[$:

$$\lim_{N \rightarrow \infty} \frac{H^s(\alpha N)}{s^N} = 1$$

Remark 7 In the current Section and in Section 3 we evaluate the expected hitting time to reach a specific vertex of \mathcal{H}_N . From a biological viewpoint this means to reach the optimal B-cell trait against the presented antigen. The single-peak landscape assumption has already been discussed in other mathematical models of GC reaction [66, 39, 38]. Looking for a perfect complementarity of the whole BCR to the target profile might not be really biologically significant : the matching of entire strings means designing a receptor for each possible antigen, this is not reasonable considering repertoire sizes. Therefore, we evaluate the hitting time of a set of vertices instead. This implies, of course, a speed-up of the time-scales (see Table 1 for instance). Let $A_r := \{\mathbf{x}_i \in \mathcal{H}_N \mid h(\mathbf{x}_i, \bar{\mathbf{x}}) \leq r\}$ be the sphere of radius r in the graph metric, centered in the target vertex $\bar{\mathbf{x}}$, and considering \mathcal{P} as transition probability matrix. We are interested in explicitly evaluate the mean hitting time to enter A_r . We consider the distances process defined in Section 2.2, hence the graph underlined by matrix \mathcal{Q} (Proposition 1). The sphere A_r can be denoted as:

$$A_r := \{j \in \{0, \dots, N\} \mid j \leq r\}$$

389 We denote by $H_i(r)$ the expected time to reach A_r starting from initial Hamming distance i . By using Equation (1), we obtain:

$$\begin{cases} H_i(r) = 0 & \text{if } i \leq r \\ H_i(r) = 1 + \frac{i}{N} H_{i-1}(r) + \frac{N-i}{N} H_{i+1}(r) & \text{if } i > r \end{cases} \quad (16)$$

Let us define $\Delta_r(i)$ as the difference between $H_i(r)$ and $H_{i-1}(r)$:

$$\Delta_r(i) := H_i(r) - H_{i-1}(r)$$

391 Therefore:

$$\begin{aligned} \Delta_r(i) &= 1 + \frac{i}{N} H_{i-1}(r) + \frac{N-i}{N} H_{i+1}(r) - H_{i-1}(r) \\ &= 1 + \frac{N-i}{N} (H_{i+1}(r) - H_{i-1}(r)) \\ &= 1 + \frac{N-i}{N} (\Delta_r(i+1) + \Delta_r(i)) \end{aligned}$$

392 And finally:

$$\Delta_r(i) = \frac{N-i}{i} \Delta_r(i+1) + \frac{N}{i} \quad (17)$$

393 With the condition:

$$\Delta_r(N) := H_N(r) - H_{N-1}(r) = 1 + H_{N-1}(r) - H_{N-1}(r) = 1 \quad (18)$$

394

395 **Theorem 6** For all $i > r \geq 0$ the mean hitting time to reach A_r starting from
396 initial Hamming distance i from $\bar{\mathbf{x}}$ is given by:

$$H_i(r) = \sum_{s=r+1}^i \frac{\sum_{j=0}^{N-s} C_N^j}{C_{N-1}^{N-s}} \quad (19)$$

Table 1: Average expected times to reach the sphere A_r of radius r centered in $\bar{\mathbf{x}}$, for different values of r . Simulations correspond to $N = 10$ and an initial Hamming distance $h(\mathbf{X}_0, \bar{\mathbf{x}}) = 10$. Table 1 shows results obtained over 20480 simulations. We denote by $|A_r|$ the number of vertices of \mathcal{H}_N included in A_r . $H_{10}(r)$ corresponds to the theoretical value obtained by Equation (19). We denote by $\widehat{\tau}_{\{\bar{\mathbf{x}}\}_n}$ the average value obtained over $n = 20480$ simulations and by $\widehat{\sigma}_n$ its corresponding estimated standard deviation.

r	$ A_r $	$H_{10}(r)$	$\widehat{\tau}_{\{\bar{\mathbf{x}}\}_n}$	$\frac{\widehat{\sigma}_n}{\sqrt{n}}$
0	1	1186.540	1184.499	8.1736
1	11	163.540	163.747	1.064
2	56	50.984	51.729	0.298
3	176	24.095	24.118	0.116

397 *Remark 8* One can demonstrate that $H_i(0) = H(i)$ as defined by Equation
398 (10).

399 *Proof* Considering Equations (17) and (18) we can demonstrate by iteration
400 that $\forall k \in \{0, \dots, N-1\}$:

$$\Delta_r(N-k) = \frac{1}{C_{N-1}^k} \sum_{j=0}^k C_N^j \quad (20)$$

401 The result follows by observing:

$$H_i(r) = \sum_{s=r+1}^i \Delta_r(s) = \sum_{s=r+1}^i \Delta_r(N - (N-s)) \quad (21)$$

□

402 We simulate the average expected time to reach a sphere of radius r centered in the vertex $\bar{\mathbf{x}}$, for different values of r . Table 1 shows the results obtained over more than 20000 simulations. We clearly see that the average hitting time decreases significantly if we consider bigger radius r , as expected.
406

407 **3 More mutational models: how does the structure of the hypercube change?**
408

409 In this section, we explore other mutation rules, which change the internal graph structure of the hypercube, therefore the dynamics of the RW and the characteristic time-scales of the exploration of the state-space.
411

412 3.1 Study of various mutation rules

413 We propose and study four mutational rules:

- 414 – a model of permutation of two bits;
- 415 – a model of switch of k -length strings;
- 416 – a model of switch of 1 or 2-length strings depending on the Hamming distance to a fixed node representing the antigen target cell;
- 417 – multiple point mutations models.

419 3.1.1 The exchange mutation model.

420 We consider a model where given an initial B-cell representing string, each mutation step consists in permuting two randomly chosen bits.
421

Definition 9 Let $\mathbf{X}_n \in \{0,1\}^N$ be the BCR at step n . Let $i \in \{1, \dots, N\}$, $j \in \{1, \dots, N\} \setminus \{i\}$ two randomly chosen indexes. We can suppose, without loss of generality, that $j > i$:

$$\mathbf{X}_{n+1} = (X_{n,1}, \dots, X_{n,i-1}, X_{n,j}, X_{n,i+1}, \dots, X_{n,j-1}, X_{n,i}, X_{n,j+1}, \dots, X_{n,N})$$

422 With this mutation rule, we loose a very important property : the connectivity of the graph. We denote by $\mathcal{H}_{(s)} \subset \{0,1\}^N$ the set containing the C_N^s vertices having s 1 in their strings. The state-space $\{0,1\}^N$ is divided into
424 $N + 1$ connected components: $\mathcal{H}_{(s)}$, $0 \leq s \leq N$.
425

426 **Proposition 6** *There are exactly $\frac{N(N-1)}{2}$ (non-oriented) edges ending at each vertex counting the possible loops. Each node $\mathbf{x} \in \mathcal{H}_{(s)}$ has exactly $\frac{(N-s)^2 - (N-s^2)}{2}$ loops.*
427
428

429 **Corollary 2** $\mathbb{P}(\mathbf{X}_n = \mathbf{x}_j | \mathbf{X}_{n-1} = \mathbf{x}_j) = \frac{(N-s)^2 - (N-s^2)}{N(N-1)}$. In particular the probability of remaining on the same node is 1 if $s = 0$ or $s = N$.
430

Proof (Proposition 6) The first statement is obtained by simple combinatory arguments. Let us consider $\mathbf{x} \in \mathcal{H}_{(s)}$ with $0 \leq s \leq N$: it is composed by exactly s ones and $N - s$ zeros. For the sake of clarity let us consider that $\{0, \dots, N\} = I \sqcup J$ so that $|I| = s$, $|J| = N - s$ and $x_i = 1 \forall i \in I$, $x_j = 0 \forall j \in J$. We obtain a loop each time we choose both random indices either in I (C_s^2 possibilities) or in J (C_{N-s}^2 possibilities). Then the total number of loops is obtained by the sum of these two cases, *i.e.* $\frac{(N-s)^2 - (N-s^2)}{2}$. \square

431 We can also describe qualitatively the behavior of the (D_n) process refer-
 432 ring to this current model. As a general principle, we have that $D_n = D_{n-1} + i$,
 433 $i \in \{0, \pm 2\}$. Therefore, clearly $\mathbb{P}(D_n = d' | D_{n-1} = d) = 0$ if $|d' - d| > 2$ or
 434 $|d' - d| = 1$. Moreover, we have maximal and minimal values of D_n depending
 435 on s_0 and \bar{s} so that $\mathbf{X}_0 \in \mathcal{H}_{(s_0)}$ and $\bar{\mathbf{x}} \in \mathcal{H}_{(\bar{s})}$. Indeed:

Proposition 7 *Given $\bar{\mathbf{x}} \in \mathcal{H}_{(\bar{s})}$ and $\mathbf{X}_0 \in \mathcal{H}_{(s_0)}$, then $\forall n \geq 0$:*

$$\begin{cases} |\bar{s} - s_0| \leq D_n \leq \bar{s} + s_0 & \text{if } \bar{s} + s_0 \leq N \\ |\bar{s} - s_0| \leq D_n \leq (N - \bar{s}) + (N - s_0) & \text{if } \bar{s} + s_0 > N \end{cases}$$

Proof The proof follows immediately by counting how many possibilities there are to arrange s ones and $N - s$ zeros in a N -length string. \square

Remark 9 From Proposition 7 one can see that if $\bar{s} = s_0 =: s$ and $2s \neq N$ then:

$$0 \leq D_n < N$$

436 *3.1.2 Class switch of k -length strings.*

437 Let $\mathbf{X}_0 = (X_{0,1}, \dots, X_{0,N}) \in \{0, 1\}^N$ be the B-cell entering the somatic hyper-
 438 mutation process. At each mutation step we switch the class of k consecutive
 439 amino-acids.

440 **Definition 10** Let $\mathbf{X}_n \in \{0, 1\}^N$ be the BCR at step n . Let $i \in \{1, \dots, N -$
 441 $(k - 1)\}$ be a randomly chosen index. Then $\mathbf{X}_{n+1} := (X_{n,1}, \dots, X_{n,i-1}, 1 -$
 442 $X_{n,i}, \dots, 1 - X_{n,i+k-1}, X_{n,i+k}, \dots, X_{n,N})$.

443 *Remark 10* If $k = 1$ we are in the case of a SRW on \mathcal{H}_N .

444 If $k = N$ we stay on a 2-length cycle. Indeed we have that $\mathbf{X}_l = \mathbf{X}_0$ for l even
 445 and $X_l = \mathbf{1} - \mathbf{X}_0$ for l odd. For this reason the case $k = N$ does not appear
 446 interesting neither from a mathematical nor from a biological point of view.

447 Here below we give some basic properties of this RW, that one can easily
 448 prove by simple combinatory arguments.

Proposition 8 *Each vertex has exactly $N - (k - 1)$ neighbors and no loops. Therefore, for all $\mathbf{x}_i, \mathbf{x}_j$ in $\{0, 1\}^N$:*

$$\mathbb{P}(\mathbf{X}_n = \mathbf{x}_j | \mathbf{X}_{n-1} = \mathbf{x}_i) =: p_k(i, j) = \begin{cases} \frac{1}{N - (k - 1)} & \text{if } \mathbf{x}_j \sim \mathbf{x}_i \\ 0 & \text{otherwise} \end{cases}$$

449 *Remark 11* As regards to this current model, given $\mathbf{x}_i, \mathbf{x}_j \in \{0, 1\}^N$, we have:
 450 $\mathbf{x}_i \sim \mathbf{x}_j \Leftrightarrow h(\mathbf{x}_i, \mathbf{x}_j) = k$ and the k different elements have consecutive indexes.

451 Thus, $\mathcal{P}_k = (p_k(\mathbf{x}_i, \mathbf{x}_j))_{\mathbf{x}_i, \mathbf{x}_j \in \mathcal{H}_k}$ is the $2^N \times 2^N$ transition probability ma-
 452 trix.

453
 454 For fixed $k \in \{1, \dots, N\}$ the graph underlying the RW corresponding to
 455 the model of class switch of k -length strings has exactly 2^{k-1} connected com-
 456 ponents, each one composed of $2^{N-(k-1)}$ elements.

457 Because of the non connectivity of the graph, we can focus on the connected
 458 component to which \mathbf{X}_0 belongs and find out the properties of our RW on
 459 it. For fixed N and k and dealing with each connected component separately,
 460 we are describing a SRW on a $(N - (k - 1))$ -hypercube. Henceforth we obtain
 461 2^{k-1} distinct hypercube-type structures of the same size.

462
 463 We can limit our study to the connected component containing \mathbf{X}_0 , which
 464 is, up to a change of variables, a $(N - (k - 1))$ -dimensional hypercube. Let
 465 $\bar{\mathcal{P}}_k$ be the restriction of \mathcal{P}_k to this connected component. If we conveniently
 466 order the $2^{N-(k-1)}$ distinct vertices, than $\bar{\mathcal{P}}_k = \mathcal{P}_{N-(k-1)}$. At this stage, it is
 467 possible to translate all classical results we know about the SRW on \mathcal{H}_n , for $n =$
 468 $N - (k - 1)$, on each connected component of this current graph, remembering
 469 the definition of neighborhood given in Remark 11.

470 *3.1.3 Class switch of 1 or 2-length strings depending on the Hamming*
 471 *distance to $\bar{\mathbf{x}}$.*

472 The models we described in Sections 3.1.1 and 3.1.2 present an important lim-
 473 itation: the underlying graphs are non-connected. Due to our choice of affinity,
 474 a model which does not enable to explore the whole state-space is not very
 475 relevant. Indeed, if the graph is non-connected and the target chain does not
 476 belong to the connected component containing the B-cell which first enters the
 477 somatic hypermutation process, then we never reach the target configuration.
 478 From a biological viewpoint, it may be more relevant to consider a smoother
 479 affinity model, in which the BCR representing string reaches the target when
 480 most, but not all, bits are similar. In this case, considering a non-connected
 481 graph, is not necessarily a problem.
 482

Another way to overcome the problem of non-connectivity is to consider a model which allows to vary the length of the strings submitted to switch-type mutations. Moreover, it is biologically credible that during the GC process B-cells can modify their mutation rate, making it somehow proportional to their affinity to the antigen [11, 8, 31]. Indeed, B-cells compete for different rescue signals (from Helper T-cells or FDCs), and that determines their fate: undergo further mutations or differentiate into plasma cells or memory cells ([1], Chapters 7). Here we suppose that the mutational rate is inversely proportional to the affinity: the greater the affinity, the lower is the mutational rate. We found the hypothesis that the regulation of the hypermutation process is dependent on receptor affinity also in other works, as [16, 2], where the authors proposed computational implementations of the clonal selection principle to design genetic optimization algorithms, taking into account AAM during an adaptive immune response. In terms of our mathematical model, we can translate it by making the size k of the strings which can mutate to be directly proportional to the Hamming distance to \bar{x} at each mutation step:

$$k_n = f(D_n), \text{ with } f : \{0, \dots, N\} \rightarrow \{0, \dots, N\} \text{ being an increasing function.}$$

483 Despite many choices of the function f are possible, hereinafter we consider a
484 very elementary case, where f is a step function on two intervals.

Definition 11 Let $\mathbf{X}_n \in \{0, 1\}^N$ be the BCR at step n . We denote by k_n :

$$k_n := f(D_n) = \begin{cases} 1 & \text{if } D_n \leq 1 \\ 2 & \text{if } D_n > 1 \end{cases}$$

485 Let $i \in \{1, \dots, N - (k_n - 1)\}$ be a randomly chosen index. Then:
486 $\mathbf{X}_{n+1} := (X_{n,1}, \dots, X_{n,i-1}, 1 - X_{n,i}, \dots, 1 - X_{n,i+k_n-1}, X_{n,i+k_n}, \dots, X_{n,N})$.

487 This model is an interesting and simple way to generalize the basic mu-
488 tational model without losing the property of connectivity of the graph. The
489 addition of this flexibility was not only motivated by biological reasons, but we
490 also expect that this modification decreases the hitting time to a fixed node.
491 This is actually true: the hitting time is halved compared to the basic model
492 (at least for N big enough). We will also show that the stationary distribution
493 is concentrated on a half part of the hypercube, the one to whom \bar{x} belongs.

494 *Remark 12* For fixed N and $k = 2$ the graph is divided into two connected
495 components composed of 2^{N-1} vertices. Two nodes belonging to the same
496 connected component have a Hamming distance of $2t$ with $0 \leq t \leq \lfloor N/2 \rfloor$. On
497 the other hand, two vertices belonging to different connected components have
498 a Hamming distance of $(2t + 1)$ with $0 \leq t \leq \lfloor (N - 1)/2 \rfloor$.

499 In order to analyze this process, we have to distinguish two cases. For fixed
500 N and \bar{x} , the process we obtain:

501 **case 1: $D_0 = 2t$, $t > 0$.** \mathbf{X}_0 belongs to the same connected component as $\bar{\mathbf{x}}$,
 502 so we are working on a $(N - 1)$ -dimensional hypercube, following the model
 503 of class switch of 2-length strings. we stay in this connected component all
 504 over the process till we arrive at $\bar{\mathbf{x}}$, as it is impossible to obtain a Hamming
 505 distance equal to 1.

506 **case 2: $D_0 = 2t + 1$, $t > 0$.** We necessarily take $k = 2$ and Remark 12 im-
 507 plies that \mathbf{X}_0 belongs to a different connected component than $\bar{\mathbf{x}}$. In order
 508 to reach the connected component containing $\bar{\mathbf{x}}$, we have to visit a node \mathbf{x}^*
 509 so that $h(\mathbf{x}^*, \bar{\mathbf{x}}) = 1$, and $|\{\mathbf{x}^* \mid h(\mathbf{x}^*, \bar{\mathbf{x}}) = 1\}| = N$. Then, if $D_0 = 1$ we are
 510 allowed to change only one element of the B-cell representing string. With
 511 probability $1/N$ we arrive directly at $\bar{\mathbf{x}}$ and with probability $(N - 1)/N$ we
 512 obtain $D_1 = 2$. Then we go back to case 1.

513 **Proposition 9** *The graph corresponding to the current model is divided into*
 514 *two connected components: $\mathcal{H}_N^{(1-2)}$ and its complementary $\overline{\mathcal{H}_N^{(1-2)}}$, s.t. $\bar{\mathbf{x}} \in$*
 515 *$\overline{\mathcal{H}_N^{(1-2)}}$. $\overline{\mathcal{H}_N^{(1-2)}}$ is accessible from $\mathcal{H}_N^{(1-2)}$, but not conversely. Vertices be-*
 516 *longing to $\overline{\mathcal{H}_N^{(1-2)}}$ are positive recurrent and vertices belonging to $\mathcal{H}_N^{(1-2)}$ are*
 517 *transient.*

Proof The existence of two connected components depends on the use of the model of switch of 2-length strings. Indeed the structure of the graph we are considering here essentially corresponds to that of the graph underlying the model of switch of 2-length strings, up to the addition of some oriented edges from $\mathcal{H}_N^{(1-2)}$ to $\overline{\mathcal{H}_N^{(1-2)}}$. As long as we stay in $\overline{\mathcal{H}_N^{(1-2)}}$ or $\mathcal{H}_N^{(1-2)}$ we are just allowed to switch 2-length strings. Moreover, we have already observed that when we are in $\overline{\mathcal{H}_N^{(1-2)}}$ we can't exit, while when we are in $\mathcal{H}_N^{(1-2)}$ we can reach $\overline{\mathcal{H}_N^{(1-2)}}$ by visiting one among the N nodes having Hamming distance 1 from $\bar{\mathbf{x}}$, and that happens in a finite number of steps. Therefore:

$$\begin{cases} \mathbb{P}(\tau_{\mathbf{x}_i} < \infty) = 1 \text{ for all } \mathbf{x}_i \in \overline{\mathcal{H}_N^{(1-2)}} \Rightarrow \mathbf{x}_i \text{ is recurrent} \\ \mathbb{P}(\tau_{\mathbf{x}_i} < \infty) < 1 \text{ for all } \mathbf{x}_i \in \mathcal{H}_N^{(1-2)} \Rightarrow \mathbf{x}_i \text{ is transient} \end{cases}$$

In particular, vertices belonging to $\overline{\mathcal{H}_N^{(1-2)}}$ are positive recurrent as the chain is irreducible on $\overline{\mathcal{H}_N^{(1-2)}}$ and $|\overline{\mathcal{H}_N^{(1-2)}}| < \infty$. \square

518 The following known result about stochastic processes, justifies Corollary
 519 3 below.

Theorem 7 *Let $(\mathbf{X}_n)_{n \geq 0}$ be a Markov chain on a state-space \mathcal{S} and $\mathbf{x}_i \in \mathcal{S}$ be positive recurrent. Let m_i be the mean return time: $m_i = \mathbb{E}(\tau_{\{\mathbf{x}_i\}} \mid \mathbf{X}_0 = \mathbf{x}_i)$. Denoting by $\mathcal{S}_r \subseteq \mathcal{S}$ the positive recurrent connected component to which \mathbf{x}_i belongs, then a stationary distribution $\bar{\pi}$ is given by:*

$$\begin{aligned} \bar{\pi}_i &= m_i \quad \forall \mathbf{x}_i \in \mathcal{S}_r \\ \bar{\pi}_i &= 0 \quad \forall \mathbf{x}_i \in \mathcal{S} \setminus \mathcal{S}_r \end{aligned}$$

520 Theorem 7 is proven by considering the relations among recurrent and
 521 transient classes, stationary distributions and return time (see [55] for some
 522 more details).

523 **Corollary 3** *The stationary distribution for the RW we describe in the present*
 524 *section, $\bar{\pi}$, is given by:*

$$\bar{\pi}_i = \begin{cases} \frac{1}{2^{N-1}} & \text{if } \mathbf{x}_i \in \overline{\mathcal{H}}_N^{(1-2)} \\ 0 & \text{if } \mathbf{x}_i \in \mathcal{H}_N^{(1-2)} \end{cases} \quad (22)$$

525 Corollary 3 is a consequence of Theorem 7 and the study of the SRW on
 526 an N -dimensional hypercube.

527 3.1.4 Allowing 1 to k mutations

528 In this section we analyze how the N -dimensional hypercube changes if we
 529 allow 1 to k independent switch-type mutations at each step, with k fixed,
 530 $k \leq N$.

531 **Definition 12** Let $\mathbf{X}_n \in \{0, 1\}^N$ be the BCR at step n . Let k be an integer,
 532 $1 \leq k \leq N$ and $\forall i, 1 \leq i \leq k, a_i := \mathbb{P}(i \text{ independent switch-type mutations})$.
 533 Then with probability a_i, \mathbf{X}_{n+1} is obtained from \mathbf{X}_n by repeating i times,
 534 independently, the process described by Definition 5.

535 By definition, the corresponding transition probability matrix is a convex
 536 combination of \mathcal{P}^i , for $1 \leq i \leq k$ (\mathcal{P}^i is the transition probability matrix
 537 corresponding to i iterations of the process of a single bit mutation):

$$\sum_{i=1}^k a_i \mathcal{P}^i, \quad \text{with } \sum_{i=1}^k a_i = 1. \quad (23)$$

538 **Definition 13** Let us fix $a_i = 1/k \forall i$. We denote by $\mathcal{P}^{(k)} := 1/k \sum_{i=1}^k \mathcal{P}^i$.
 539 Accordingly, we denote the graph underlying this RW $\mathcal{H}_N^{(k)}$.

540 *Remark 13* Since the mutations are assumed to be independent, then k represents the maximum Hamming distance the process can cover in a single mutation step. Thanks to the independence of each single mutation, two or more mutations may nullify their respective action: in particular for $k \geq 2$ there is a non-zero probability of remaining at the same position. From a biological point of view, this behavior can be interpreted as the possibility of doing mutations which have no effect on the BCR structure.

547 We can now evaluate the eigenvalues of $\mathcal{P}^{(k)}$, $\lambda_j^{(k)}$ by using the eigenvalues
 548 λ_j of \mathcal{P} (Section 2.1). Due to the fact that all \mathcal{P}^i commute with each other,
 549 the eigenvalues are given by:

$$\lambda_j^{(k)} = \frac{1}{k} \sum_{i=1}^k \lambda_j^i \quad (24)$$

550 and $\mathcal{P}^{(k)}$ and \mathcal{P} have the same eigenvectors. We give explicitly the expression
 551 of all $\lambda_i^{(k)}$ and concentrate on the second largest eigenvalue, $\lambda_2^{(k)}$.

552 **Proposition 10** *The $N+1$ distinct eigenvalues of matrix $\mathcal{P}^{(k)}$ are:*

$$\begin{aligned} 553 & - \lambda_1^{(k)} = 1 ; \\ 554 & - \lambda_j^{(k)} = \frac{\lambda_j}{k} \cdot \frac{1 - \lambda_j^k}{1 - \lambda_j} \text{ for } 2 \leq j \leq N ; \\ 555 & - \lambda_{N+1}^{(k)} = \frac{1}{2k} \left((-1)^k - 1 \right) = \begin{cases} 0 & \text{if } k \text{ is even} \\ -1/k & \text{if } k \text{ is odd} \end{cases} \end{aligned}$$

556 *The multiplicity of $\lambda_j^{(k)}$ is $\binom{N}{j-1}$, $1 \leq j \leq N+1$*

Proof This result comes directly from the evaluation of Equation (24), for the already known values of all λ_j (Corollary 1). \square

557 Then, in particular, the second largest eigenvalue of $\mathcal{P}^{(k)}$ is:

$$\lambda_2^{(k)} = \frac{N-2}{2k} \left(1 - \left(1 - \frac{2}{N} \right)^k \right) \quad (25)$$

Remark 14 For all $k \geq 2$, $\lambda_2 > \lambda_2^{(k)}$. First of all, we can observe that $\lambda_2^{(k)}$ decreases for increasing k . Therefore:

$$\lambda_2 - \lambda_2^{(k)} \geq \lambda_2 - \lambda_2^{(2)} = \frac{N-2}{4N^2} (4N - N^2 + (N-2)^2) = \frac{N-2}{N^2} > 0$$

558 For $N \gg 1$, the series expansion of $\lambda_2^{(k)}$ gives us:

$$\begin{aligned} \lambda_2^{(k)} &= \frac{N-2}{2k} \left(1 - \left(1 - \frac{2k}{N} + \frac{2k(k-1)}{N^2} + \mathcal{O}\left(\frac{1}{N^3}\right) \right) \right) \\ &= \frac{N-2}{N} - \frac{(N-2)(k-1)}{N^2} + \mathcal{O}\left(\frac{1}{N^2}\right) \end{aligned}$$

559 We can observe how the spectral gap changes. If we consider the series
 560 expansion of $\left(1 - \frac{2}{N}\right)^k$ for $N \rightarrow \infty$, we get:

$$\lambda_1^{(k)} - \lambda_2^{(k)} = \frac{2}{N} + \frac{(N-2)(k-1)}{N^2} + \mathcal{O}\left(\frac{1}{N^2}\right)$$

561 It can be interesting to choose k as a function of N . Let us consider, for
 562 example, $k = \alpha N$, with $0 < \alpha \leq 1$. In this case, we have:

$$\begin{aligned} \lambda_2^{(\alpha N)} &= \frac{N-2}{2\alpha N} \left(1 - \left(1 - \frac{2}{N} \right)^{\alpha N} \right) \\ &\stackrel{\text{for } N \rightarrow \infty}{=} \frac{N-2}{2\alpha N} \left(1 - \left(e^{-2\alpha} + \mathcal{O}\left(\frac{1}{N}\right) \right) \right) \\ &= \frac{(N-2)(1-e^{-2\alpha})}{2\alpha N} + \mathcal{O}\left(\frac{1}{N}\right) \rightarrow \frac{1-e^{-2\alpha}}{2\alpha} \text{ for } N \rightarrow \infty \end{aligned}$$

563 We can observe that $\frac{1-e^{-2\alpha}}{2\alpha} =: \bar{\lambda}_2^{(\alpha N)}$ decreases when α increases. More-
 564 over:

- 565 - $\bar{\lambda}_2^{(\alpha N)} \rightarrow 1$ for $\alpha \rightarrow 0$, which means that the spectral gap, $1 - \lambda_2^{(\alpha N)}$ con-
 566 verges to zero for $N \rightarrow \infty$ and $\alpha \rightarrow 0$;
- 567 - If $\alpha = 1$ then $\bar{\lambda}_2^{(N)} = \frac{1}{2} - \frac{1}{2e^2}$. Therefore, the spectral gap is $\frac{1}{2} + \frac{1}{2e^2}$

568 The spectral gap indicates how quickly a RW converges to its stationary
 569 distribution. As expected, if $\alpha \rightarrow 0$ then the spectral gap gets close to 0. On the
 570 other hand for all $\alpha > 0$ the spectral gap tends to a strictly positive quantity,
 571 while the spectral gap corresponding to the case of the basic model converges
 572 to zero for $N \rightarrow \infty$. In particular, when $\alpha = 1$ (*i.e.* we are considering the
 573 optimal case, in which we are allowed to do among 1 and N mutations at each
 574 mutation step), the spectral gap, $\frac{1}{2} + \frac{1}{2e^2}$, is significantly bigger than the one
 575 obtained for the basic model, $2/N$.

576 3.2 Comparison of hitting times

577 In this section we compare hitting times referring to some relevant mutational
 578 models we have already presented. We do not consider models that entail
 579 non-connected graphs (the model of class switch of k -length strings and the
 580 exchange mutation model): this choice is motivated by the discussion from the
 581 beginning of Section 3.1.3. In Table 2 we collect most important characteristics
 582 of these RWs on $\{0,1\}^N$: the hitting time and its approximation for big N ,
 583 that we will discuss in this current section, the stationary distribution and the
 584 value of the second larger eigenvalue when known.

585 3.2.1 Class switch of 1 or 2-length strings depending on the Hamming 586 distance to $\bar{\mathbf{x}}$.

587 We use results obtained in Section 2 for the (D_n) process concerning the
 588 SRW on the N -dimensional hypercube and we apply them to this model.
 589 Here we shall introduce another definition of the distance, which is adapted
 590 to a connected component $\mathcal{H}_{N,2} \subset \{0,1\}^N$, where $\mathcal{H}_{N,2}$ denotes one of the

Table 2: Table 2 summarizes the main characteristics of most random processes we introduce and analyze in Sections 2 and 3.

Model	Hitting time	Stationary distribution	Second biggest eigenvalue
Basic model	$H(\bar{d}) = \sum_{d=0}^{\bar{d}-1} \frac{\sum_{j=1}^{N-1-d} C_N^{d+j+1}}{C_N^d} \sim 2^N$	π	$1 - \frac{2}{N}$
Switch 1-2	$\sim 2^{N-1}$	$\pi _{\overline{\mathcal{H}_N}^{(1-2)}}$	-
Allowing 1 to k mutations	$\overline{T}_N^{(k)}(\bar{d}) = \sum_{l=2}^{2^N} \mu_l^{(k)} R_N(l, \bar{d})$ $\frac{1}{2^N C_N^{\bar{d}}} \sum_{l=2}^{2^N} \mu_l^{(k)} R_N(l, \bar{d})$	π	$\frac{N-2}{2k} \left(1 - \left(\frac{N-2}{N}\right)^k\right)$

591 two parts in which $\{0,1\}^N$ is divided applying the model of class switch of
592 2-length strings (Section 3.1.2). We recall that $\mathcal{H}_{N,2}$ is a $(N-1)$ -dimensional
593 hypercube, and that the graph underlying the model of class switch of 1 or
594 2-length strings corresponds essentially to the graph obtained with the model
595 of switch of 2-length strings, up to the addition of some oriented edges from
596 $\mathcal{H}_N^{(1-2)}$ to $\overline{\mathcal{H}_N}^{(1-2)}$.

597 **Definition 14** For all $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{H}_{N,2}$ we denote by $h^{(2)}(\mathbf{x}_i, \mathbf{x}_j)$ the number
598 of edges in a shortest path connecting them. Simultaneously we denote by
599 $D_n^{(2)} = h^{(2)}(\mathbf{X}_n, \bar{\mathbf{x}})$, $D_n^{(2)} \in \{0, \dots, N-1\} \forall n \geq 0$.

600 Considering the process $(D_n^{(2)})_{n \geq 0}$, all results stated in Section 2 hold
601 true. Furthermore, let us denote by $\mathbb{E}_{\mathbf{x}_i}^{(2)}[\tau_A]$ the expected number of steps
602 before set $A \in \mathcal{H}_{N,2}$ is visited starting at $\mathbf{x}_i \in \mathcal{H}_{N,2}$ and following the model
603 of switch of 2-length strings. Then, we also denote by $H_{N-1}^{(2)}(d) = \mathbb{E}_{\mathbf{x}}^{(2)}[\tau_{\{\bar{\mathbf{x}}\}}]$
604 where $d = h^{(2)}(\mathbf{x}, \bar{\mathbf{x}})$.

605 *Remark 15* Clearly if $D_0 = 2t$ and $t > 0$, which means that \mathbf{X}_0 and $\bar{\mathbf{x}}$ belong to
606 the same connected component in the model of class switch of 2-length strings,
607 then the mean hitting time for the current model will be of the order of a half
608 the mean hitting time for the basic model. Indeed, since we are considering
609 here a $(N-1)$ -dimensional hypercube instead of a N -dimensional one.

610 The result below, which is an immediate application of the Ergodic Theo-
611 rem, will help us understand better the general behavior of this mean hitting
612 time:

613 **Proposition 11** Let $(\mathbf{X}_n)_{n \geq 0}$ be a SRW on \mathcal{H}_N . We denote by $T_d^+ := \inf\{n \geq$
614 $1 \mid D_n = d\}$ and $T_d := \inf\{n \geq 0 \mid D_n = d\}$. Then:

$$\mathbb{E}_{D_0=d}[T_d^+] = \frac{2^N}{C_N^d} \quad (26)$$

Proof The proof is obtained by applying the Ergodic Theorem to the (D_n) process and its stationary distribution, the binomial probability distribution. \square

615 For the discussion we made in Section 2.2 and, in particular, Remark 3 we
616 can conclude that for $N \gg 1$ the order of magnitude of the time we spend to
617 reach the N nodes at Hamming distance 1 from $\bar{\mathbf{x}}$ is:

$$\mathbb{E}_{D_0=d}[T_1] \sim \frac{2^N}{N} \quad (27)$$

618 Then we can claim the following result, which comes directly from Equation
619 (27):

Proposition 12 *Let us suppose that $D_0 = 2t^* + 1$ with $0 < t^* \leq \lfloor (N-1)/2 \rfloor$. Then for $N \gg 1$ we have:*

$$\mathbb{E}_{D_0=d}^{(2)}[T_1] \sim \frac{2^{N-1}}{N}$$

620 Finally:

Proposition 13 *We denote by $\mathbb{E}_{\mathbf{x}_0}^{(1-2)}[\tau_{\{\bar{\mathbf{x}}\}}]$ the mean hitting time to reach $\bar{\mathbf{x}}$ starting from \mathbf{x}_0 and referring to the mutation model of class switch of 1 or 2 length strings. Then, for $N \gg 1$ we have:*

$$\mathbb{E}_{\mathbf{x}_0}^{(1-2)}[\tau_{\{\bar{\mathbf{x}}\}}] \sim \frac{1}{2} \mathbb{E}_{\mathbf{x}_0}[\tau_{\{\bar{\mathbf{x}}\}}] \quad \text{with} \quad \mathbb{E}_{\mathbf{x}_0}[\tau_{\{\bar{\mathbf{x}}\}}] \sim 2^N,$$

621 where $\mathbb{E}_{\mathbf{x}_0}[\tau_{\{\bar{\mathbf{x}}\}}]$ is the hitting time from \mathbf{x}_0 to $\bar{\mathbf{x}}$ according to the basic model,
622 as defined in Section 2.3.

Proof First of all we observe that the last statement is a direct consequence of Proposition 3. As far as the first statement is concerned, we observe that according to the model we are analyzing here and due to Proposition 12, for $N \gg 1$ the order of magnitude of $\mathbb{E}_{\mathbf{x}_0}^{(1-2)}[\tau_{\{\bar{\mathbf{x}}\}}]$ is:

$$\mathbb{E}_{\mathbf{x}_0}^{(1-2)}[\tau_{\{\bar{\mathbf{x}}\}}] \sim \frac{1}{2} \left(\frac{2^{N-1}}{N} + 2^{N-1} \right) + \frac{1}{2} 2^{N-1}$$

where the first term corresponds to the case $\mathbf{x}_0 \notin \overline{\mathcal{H}_N}^{(1-2)}$ and the second one corresponds to the opposite case (as we choose randomly the first vertex, \mathbf{x}_0 , we have probability 1/2 that it belongs to each part of the hypercube). For the last term we used again Proposition 3 applied to a $(N-1)$ -dimensional hypercube and according to the $(D_n^{(2)})$ process and the corresponding hitting time $H_{N-1}^{(2)}(d)$. The result follows. \square

Table 3: Average expected times from $[0, \dots, 0]$ to $[1, \dots, 1]$, comparing the basic mutational model and the model of class switch of 1 or 2 length strings. Here we denote by $\widehat{\tau}_{\{\bar{\mathbf{x}}\}_n}$ the average value obtained over n simulations and by $\widehat{\sigma}_n$ its corresponding estimated standard deviation.

Mutational model	N	n	$\widehat{\tau}_{\{\bar{\mathbf{x}}\}_n}$	$\frac{\widehat{\sigma}_n}{\sqrt{n}}$
Basic	10	5000	1188.7996	16.2930
	11	5000	2312.5648	32.1073
Switch 1-2	10	5000	602.8124	8.4773
	11	5000	1181.5174	16.9023

623 *Remark 16* We simulated the basic mutational model and the model of class
 624 switch of 1 or 2 length strings in order to compare the hitting times from
 625 $\mathbf{x}_0 := [0, \dots, 0]$ to $\bar{\mathbf{x}} := [1, \dots, 1]$ for both mutational models. We consider the
 626 case $N = 10$ and $N = 11$ in order to have an example in which the process
 627 starts from $\overline{\mathcal{H}}_N^{(1-2)}$ and from $\mathcal{H}_N^{(1-2)}$ respectively. Indeed, if $N = 10$ the
 628 process starts from the connected component to which $\bar{\mathbf{x}}$ belongs, while when
 629 $N = 11$ we have to reach one of the N nodes having distance 1 from $\bar{\mathbf{x}}$ to reach
 630 the connected component containing $\bar{\mathbf{x}}$. The average resulting hitting times
 631 are summarized in Table 3.

632 3.2.2 Allowing 1 to k mutations.

In this section we study the mean hitting time to cover a fixed Hamming distance d . First of all, we give the expression of the hitting time from node i to node j using the spectra. This formula is deduced by the more general one given in [49], in the case of regular graphs (the graph obtained by a convex combination of matrices \mathcal{P}^i is a regular multigraph). We refer to the notations given in Section 2 for the eigenvectors of matrix \mathcal{P} : $\mathbf{v}_s = (v_{s1}, \dots, v_{s2N})$ is the normalized eigenvector of \mathcal{P} corresponding to λ_s . These eigenvectors are the columns of matrix Q_N (Section 2.1), and each component v_{si} corresponds to node i , as they were organized while constructing the adjacency matrix. Denoting by $T(i, j)$ the hitting time from node i to node j in $\mathcal{H}_N^{(k)}$, we obtain the following expression:

$$T(i, j) = 2^N \sum_{l=2}^{2^N} \frac{1}{1 - \lambda_l^{(k)}} (v_{lj}^2 - v_{li}v_{lj}),$$

which can be written using column vectors of \mathcal{Z}_N .

$$T(i, j) = \sum_{l=2}^{2^N} \frac{1}{1 - \lambda_l^{(k)}} (z_{lj}^2 - z_{li}z_{lj})$$

633 We are interested in studying the equation below:

$$\bar{T}_N^{(k)}(d) := \frac{1}{2^N C_N^d} \sum_{h(i,j)=d} T(i,j) = \frac{1}{2^N C_N^d} \sum_{l=2}^{2^N} \frac{1}{1 - \lambda_l^{(k)}} \sum_{h(i,j)=d} (z_{l_j}^2 - z_{l_i} z_{l_j}), \quad (28)$$

634 where $2^N C_N^d$ corresponds to the number of couples of nodes of $\{0,1\}^N$ having
635 Hamming distance d .

636

637 First of all we can observe that for all l and for all j , $z_{l_j}^2 = 1$. Moreover,
638 in order to simplify notations, we denote $\mu_l^{(k)} := (1 - \lambda_l^{(k)})^{-1}$. Also, we denote
639 $R_N(l, d) := \sum_{h(i,j)=d} z_{l_i} z_{l_j}$. Finally we obtain:

Proposition 14

$$\bar{T}_N^{(k)}(d) = \sum_{l=2}^{2^N} \mu_l^{(k)} - \frac{1}{2^N C_N^d} \sum_{l=2}^{2^N} \mu_l^{(k)} R_N(l, d) \quad (29)$$

640 All the elements of this equation are known, except $R_N(l, d)$. Let us consider
641 the $2^N \times (N+1)$ matrix $\mathcal{R}_N = (R_N(l, d))$, with $1 \leq l \leq 2^N$ and $0 \leq d \leq N$. One
642 can prove by iteration:

Proposition 15

$$\mathcal{R}_N = \mathcal{Z}_N \cdot \mathcal{L}_N \quad (30)$$

where $\mathcal{Z}_N := (\mathbf{z}_1, \dots, \mathbf{z}_{2^N})$ is recursively obtained from \mathcal{Z}_{N-1} (Section 2.1),
and

$$\left\{ \begin{array}{l} \mathcal{L}_1 = 2I_2, \text{ } I_n \text{ being the } n\text{-dimensional identity matrix} \\ \mathcal{L}_N = \begin{pmatrix} 2 \cdot \mathcal{L}_{N-1} & \mathbf{0}_{2^{N-1}} \\ \mathbf{0}_{2^{N-1}} & 2 \cdot \mathcal{L}_{N-1} \end{pmatrix}, \text{ } \mathbf{0}_n \text{ being the } n\text{-length zero column vector} \end{array} \right.$$

643 3.2.3 Numerical simulations

644 In Figure 3 we plot some examples of the dependence of $\bar{T}_N^{(k)}(d)$ on d and k
645 for different values of N .

646

647 Figure 3 (a) shows that for increasing k , $\bar{T}_N^{(k)}(d)$ varies on a smaller interval:
648 $[1023, 1186.5]$ for $k = 1$, $[1028.1, 1068.6]$ for $k = 5$ and $[1025.6, 1044.8]$ for $k = 10$.
649 It is intuitive to understand this fact: the hitting time depends less from the
650 initial Hamming distance if we allow more mutations at the same mutational
651 step. Indeed, we can actually visit more distant nodes since the first steps, so
652 the initial Hamming distance has a smaller influence on the result. Figures 3

653 (b) and 3 (c) show the dependence of $\bar{T}_N^{(k)}(d)$ on k . We obtain the best result
 654 for the biggest k , except in the case $d = 1$ (as already shown by Figure 3 (a)).
 655 Curves corresponding to the case $d = 5$ and $d = 10$ are really close: we can eval-
 656 uate their minimal and maximal values, which are respectively 1043.25 and
 657 1177.60 for $d = 5$; 1044.82 and 1186.54 for $d = 10$. This fact highlights once
 658 again that if $d > 1$, the initial Hamming distance poorly influences the value
 659 of the hitting time. The case $d = 1$ shows surprisingly that the hitting time is
 660 not necessarily a monotone function of k . Figure 3 (c) allows to focus to this
 661 behavior and better understand its causes. Indeed, as N is quite small, this
 662 figure shows more clearly the oscillating behavior of $\bar{T}_N^{(k)}(d)$ while studying its
 663 dependence on k : for even values of k , $\bar{T}_5^{(k)}(1)$ increases, while for odd values
 664 of k it decreases. Intuitively, as the distance we want to cover is $d = 1$, if we
 665 allow to do 2 mutations instead of simply one, then we have a high probability
 666 to go further since the beginning of the process. Let us now look to Equation
 667 (28) and, in particular to the factor: $\sum_{l=2}^{2^N} (1 - \lambda_l^{(k)})^{-1}$. We can understand
 668 the phenomenon plotted in Figure 3 (c) by looking at Proposition 10. If k
 669 is odd and little enough then the last eigenvalue, which is negative (equal to
 670 $-1/k$), has an important negative influence over the value of $\bar{T}_N^{(k)}(d)$. Clearly,
 671 this fact has a substantial effect only if N and k are little enough, otherwise
 672 it will be compensated by the effect of all other eigenvalues.

673
 674 One may wonder what would be the best choice for the coefficients a_i (De-
 675 finition 12), $1 \leq i \leq k$, so that $\bar{T}_N^{(k)}(d)$ is minimized for a fixed k . We have
 676 to minimize the convex combination $\sum_{i=1}^k a_i \lambda_i^i$. The answer is quite evident:
 677 if $k > 2$ the minimum is obtained by taking all $a_i = 0$ and $a_{k^*} = 1$, where
 678 $k^* = 2\lfloor(k+1)/2\rfloor - 1$. Consequently, the best choice for the transition probabil-
 679 ity matrix is \mathcal{P}^{k^*} . The fact that we need to consider the greater odd component
 680 has also a more intuitive explanation. Indeed if we consider the RW given by
 681 \mathcal{P}^{2t} , we will be trapped in one of the connected-components of the graph due
 682 to the bipartite structure of the hypercube. One can remark that the graph
 683 corresponding to \mathcal{P}^{2t} is non-connected $\forall t > 0$. Therefore, we will not be able to
 684 reach those nodes having a different parity of 1s in their string, referring to \mathbf{X}_0 .

685
 686 In Figures 3 (d), 3 (e), 3 (f) and 3 (g) we plotted together the values of
 687 hitting times to cover a Hamming distance d for different values of N , k , and
 688 d , comparing the process given by $\mathcal{P}^{(k)}$ and the one corresponding to \mathcal{P}^{k^*} .
 689 This gives more evidence of the fact that the second one is the optimal one.
 690 It is interesting to look at the case in which d is fixed and we let k vary.
 691 For $k = 1$ both processes gave the same result as $\mathcal{P}^{1^*} = \mathcal{P} = \mathcal{P}^{(1)}$. Moreover,
 692 for $k = 2$ the process $\mathcal{P}^{(2)}$ is clearly the faster one: we recall that defining
 693 \mathcal{P}^{k^*} we consider the greater odd k , and then $\mathcal{P}^{2^*} = \mathcal{P}$, while the process $\mathcal{P}^{(2)}$
 694 allows to do 1 or 2 mutations at each mutation step. Then \mathcal{P}^{k^*} is the best
 695 choice among all possible convex combinations of \mathcal{P}^i iff $k > 2$. In Figures 3
 696 (d) and 3 (e) we observe the oscillating behavior of $\bar{T}_N^{k^*}(d)$. That depends

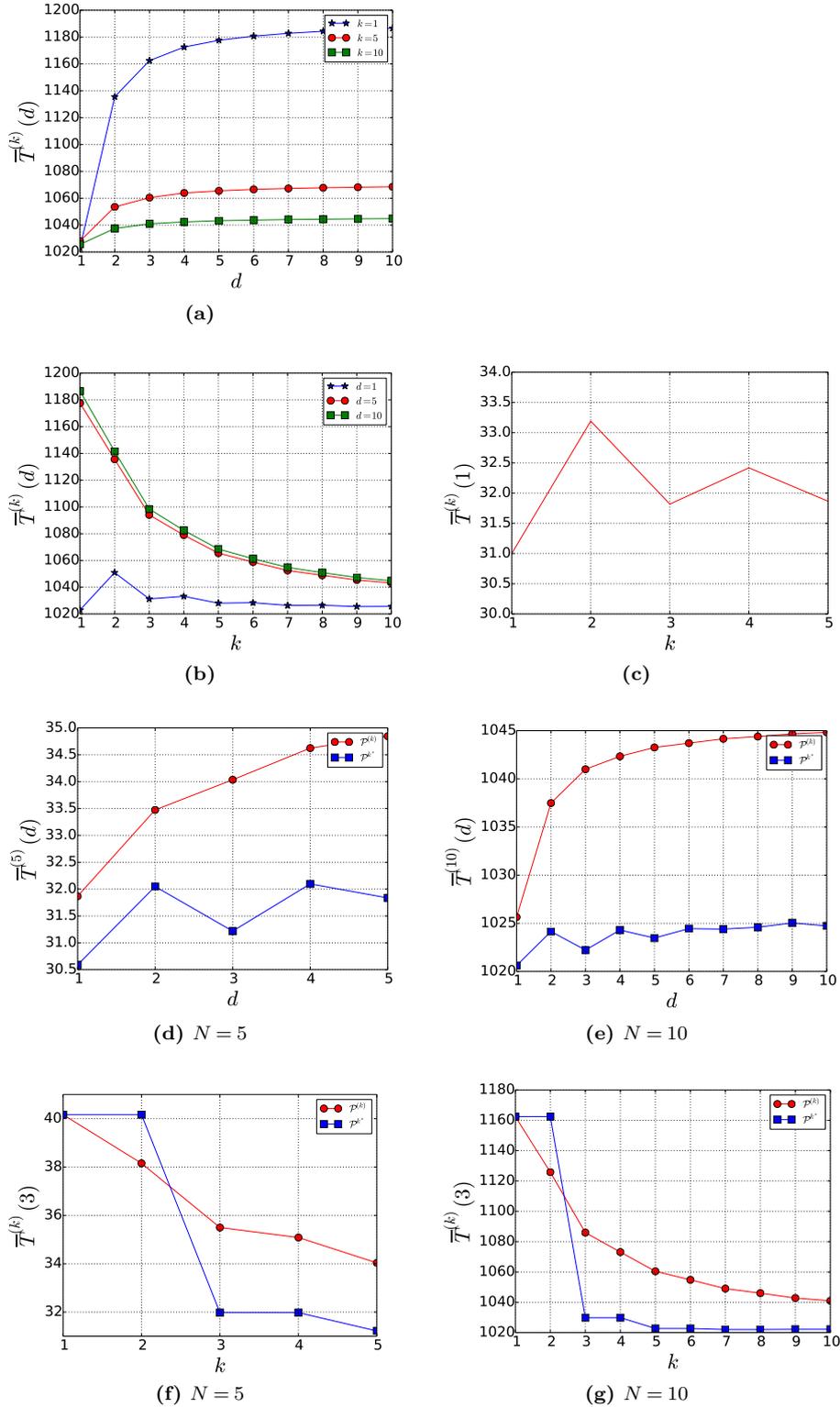


Figure 3: (a) Dependence of $\bar{T}_N^{(k)}(d)$ on d for $N = 10$ and $k = 1, 5$ or 10 . (b) Dependence of $\bar{T}_N^{(k)}(d)$ on k for $N = 10$ and different values of d . (c) Dependence of $\bar{T}_N^{(k)}(1)$ on k . (d, e) Dependence of $\bar{T}_N^{(k)}(d)$ on d for different values of both N and k . Values obtained by using as transition probability matrices $\mathcal{P}^{(k)}$ and \mathcal{P}^{k*} respectively are compared. (f, g) Dependence of $\bar{T}_N^{(k)}(d)$ on k for different values of both N and d . Again, cases corresponding to $\mathcal{P}^{(k)}$ and \mathcal{P}^{k*} are compared.

697 on the structure of \mathcal{R}_N , considering that $\sum_{l=2}^{2^N-1} R_N(l, d) = 0$ for d odd and
 698 $\sum_{l=2}^{2^N-1} R_N(l, d) = -2(2^N C_N^d)$ for d even. One can get convinced of this fact by
 699 explicitly computing $\overline{T}_N^{k^*}(d)$ for $N = 3$. Moreover simulations show that this
 700 behavior is softened for increasing d , and that $\overline{T}_N^{k^*}(N-1) > \overline{T}_N^{k^*}(N)$. This fact
 701 is confirmed by simulations on the real process. Finally, Figures 3 (f) and 3 (g)
 702 clearly show that for $k = 2$ the process given by $\mathcal{P}^{(k)}$ allows to cover quickly
 703 a fixed Hamming distance. As expected, the best hitting time is obtained for
 704 $k = N$, and for increasing N and k the value of this hitting time has a smaller
 705 variation.
 706

Table 4: An example of comparison between the theoretical and experimental values of $\overline{T}_5^{(5)}(4)$ for $\mathcal{P}^{(5)}$. $\widehat{\overline{T}_5^{(5)}(4)}_n$ denotes the average value obtained over n simulations and $\widehat{\sigma}_n$ its corresponding estimated standard deviation.

Transition probability matrix	N	d	k	n	$\overline{T}_5^{(5)}(4)$	$\widehat{\overline{T}_5^{(5)}(4)}_n$	$\widehat{\frac{\sigma_n}{\sqrt{n}}}$
$\mathcal{P}^{(k)}$	5	4	5	480000	34.62	34.67	0.05

707 We can test all these observations by simulating the real process for both
 708 transition probability matrices, \mathcal{P}^{k^*} and $\mathcal{P}^{(k)}$. Results obtained are consistent
 709 with our theoretical analysis. In order to give an idea of experimental values
 710 obtained by testing the process, in Table 4 we compare the theoretical value of
 711 $\overline{T}_N^{(k)}(d)$ corresponding to $\mathcal{P}^{(k)}$, and the experimental value with its precision,
 712 for $N = 5$, $k = 5$ and $d = 4$.

713 **4 Modeling issues**

714 The mathematical framework described in previous sections can be used to
 715 model mutations characteristic of SHM. In Sections 4.1 and 4.2 we give some
 716 more details about GCs and the binding between B-cells and antigens. There-
 717 fore, in Section 4.3 we set the modeling assumptions which justify to math-
 718 ematically describe SHMs as RWs on binary strings. Of course, this is a not
 719 exhaustive approximation. Hence, some limitations are discussed in Section
 720 4.4 and some propositions for further developments are given as well.

721 **4.1 The germinal center reaction**

722 Antigen-activated B-cells, together with their associated T cells, move into a
 723 primary lymphoid follicle, where they proliferate and ultimately form a GC.

GCs are composed mainly of B-cells, but antigen specific T-cells, which have also been activated and migrated to the lymphoid follicle, make up about 10% of GC lymphocytes and provide indispensable help to B-cells [60, 68, 54]. Indeed, when B-cells start to proliferate in GC, they need to receive proper survival signals, or they die by apoptosis. The number of B-cells within a germinal center grows at high pace: it can double every 6-8 hours [31, 19]. After about 3 days of strong proliferation, B-cells start undergoing SHM, in order to diversify the variable region of their BCRs, and those cells that express newly generated BCRs are selected for enhanced antigen binding. The fast proliferation rate of B-cells is required for the generation of a large number of modified BCRs within a short frame time (one cell gives 10^4 blasts in 72 hours). Some B-cells positively selected in the light zone differentiate into memory B-cells or plasma cells. The GC reaches its maximal size within approximately two weeks, after which the structure slowly involutes and disappears within several weeks [75]. During the GC process B-cells are subjected to powerful selection mechanisms that facilitate the generation of high affinity antibodies: a B-cell that express a newly generated BCR needs to be tested for enhanced antigen binding. This process is mediated by FDCs and follicular helper T-cells. BCR stimulation through antigen binding coupled with co-stimulatory signals transmitted by GC T-cells, provides survival signals to the cell. By contrast, failure of the BCR to bind antigen and receive proper rescue signals causes cell death by apoptosis [19]. The final differentiation of a GC B-cell into a plasma cell or a long-lived memory B-cell is driven by the acquisition of a high-affinity BCR. For short-lived memory B-cells, the differentiation process seems to be stochastic, as throughout GC reaction B-cells are constantly selected to enter the memory pool [54, 70].

4.2 B-cell receptors and antigen-antibody binding

Immunoglobulins (Ig) present at the antigen receptor are Y-shaped macro proteins composed of four polypeptide chains assembled by disulfide bonds: two identical heavy (H) chains and two identical light (L) chains. Each chain consists of two regions: a constant (C) region, which has an effector function, and a variable (V) region composed by the variable parts of the two chains together. During GC reaction the only one involved in SHMs is the V region, which also determines the antigen binding site ([54], Chapter 1). We call *antigen binding site* or *paratope* the specialized portion of the BCR V region used for identifying other molecules, while the regions on any molecule that paratopes can recognize are called *epitopes*. B-cells are able to bind ligands whose surfaces are ‘complementary’ to that of their antigen binding site, where complementarity means that the amino-acids composing the paratope and the epitope are distributed in such a way to form bonds which hold the antigen to the B-cell. In this case these bonds are all non-covalent (as hydrogen bonds, electrostatic bonds, van der Waals forces and hydrophobic bonds), which are by their nature reversible. Multiple bonding between the antigen and the B-cell

767 ensures that the antigen is bound tightly to the B-cell. The interaction between
768 paratope and epitope can be characterized in terms of a binding affinity, pro-
769 portional to their complementarity. The *affinity* is the strength of the reaction
770 between a single antigenic determinant and a single combining site on the B-
771 cell: it summarizes the attractive and repulsive forces operating between the
772 antigenic determinant and the combining site of the B-cell, and corresponds
773 to the equilibrium constant that describes the antigen-B-cell reaction [1, 78, 46].
774

775 Each antigen typically has several epitopes, so that the surface of an antigen
776 presents variable motifs that B-cells, through their receptors, can discriminate
777 as distinct epitopes. If we define an epitope by its spatial contact with a BCR
778 during binding, the number of relevant amino-acids is approximately 15, and
779 among these amino-acids only around 5 in each epitope strongly influence the
780 binding. These strong sites may contribute about one-half of the total free en-
781 ergy of the reaction, while the other amino-acids influence in binding constant
782 by up to one order of magnitude or even have no detectable effect. Simulta-
783 neously, a BCR contains a variety of possible binding sites and each antibody
784 binding site defines a paratope: about 50 variable amino-acids make up the
785 potential binding area of a BCR. In agreement with the above, only around
786 15 among these 50 amino-acids physically contact a particular epitope: these
787 define the structural paratope. Consequently, antibodies have a large num-
788 ber of potential paratopes as the 50 or so variable amino-acids composing the
789 binding region define many putative groups of 15 amino-acids [46].
790

791 Substitutions both in and away from the binding site can change the spatial
792 conformation of the binding region and affect the binding reaction. The con-
793 sequence of mutation at a particular site depends on the original amino-acid
794 and the amino-acid used for substitution ([1], Chapter 4).

795 4.3 From DNA to amino-acids: choosing the best viewpoint

796 Mutations observed on the binding site of B-cells during the GC process are
797 the result of genetic mutations produced by SHM on the portion of DNA en-
798 coding for the BCR V region. In the current section we discuss a model of
799 genetic mutations and its effects on the amino-acid string, under the assump-
800 tion of having two amino-acid classes. We show that the framework we set up
801 in previous sections can adapt to model the effects of SHM over BCRs and
802 study the variation of the affinity with the presented antigen.
803

804 The genetic code is a sequence of four nucleotides, guanine (G), adenine
805 (A) (called purines), thymine (T) and cytosine (C) (pyrimidines), joined to-
806 gether. They make three-letter words: the codons. Each codon corresponds to
807 a specific amino-acid or to a stop signal, which interrupts the building of the
808 protein during translation. As the number of possible combinations of 4 nu-
809 cleotides in 3-length words is 64, and there exists 20 amino-acids in naturally

810 derived proteins, more than a single codon codes for the same amino-acid [69].
 811 Table 5 shows the correspondence between codons and amino-acids.

812

Table 5: The correlation between codons and amino-acids: most of the amino-acids derives from more than a single codon.

	T		C		A		G		
T	TTT	Phe (F)	TCT	Ser (S)	TAT	Tyr (Y)	TGT	Cys (C)	T
	TTC	Phe (F)	TCC	Ser (S)	TAC	Tyr (Y)	TGC	Cys (C)	C
	TTA	Leu (L)	TCA	Ser (S)	TAA	Stop	TGA	Stop	A
	TTG	Leu (L)	TCG	Ser (S)	TAG	Stop	TGG	Trp (W)	G
C	CTT	Leu (L)	CCT	Pro (P)	CAT	His (H)	CGT	Arg (R)	T
	CTC	Leu (L)	CCC	Pro (P)	CAC	His (H)	CGC	Arg (R)	C
	CTA	Leu (L)	CCA	Pro (P)	CAA	Gln (Q)	CGA	Arg (R)	A
	CTG	Leu (L)	CCG	Pro (P)	CAG	Gln (Q)	CGG	Arg (R)	G
A	ATT	Ile (I)	ACT	Thr (T)	AAT	Asn (N)	AGT	Ser (S)	T
	ATC	Ile (I)	ACC	Thr (T)	AAC	Asn (N)	AGC	Ser (S)	C
	ATA	Ile (I)	ACA	Thr (T)	AAA	Lys (K)	AGA	Arg (R)	A
	ATG	Met (M)	ACG	Thr (T)	AAG	Lys (K)	AGG	Arg (R)	G
G	GTT	Val (V)	GCT	Ala (A)	GAT	Asp (D)	GGT	Gly (G)	T
	GTC	Val (V)	GCC	Ala (A)	GAC	Asp (D)	GGC	Gly (G)	C
	GTA	Val (V)	GCA	Ala (A)	GAA	Glu (E)	GGA	Gly (G)	A
	GTG	Val (V)	GCG	Ala (A)	GAG	Glu (E)	GGG	Gly (G)	G

813 Different kind of genetic mutations can affect the DNA sequence of a gene.
 814 They can be regrouped in three main categories: base substitutions, inser-
 815 tions and deletions. A single base substitution is a switch of a nucleotide with
 816 another. This is the simplest kind of mutation and it can turn out to be mis-
 817 sense, nonsense or silent, once we observe the resulting new protein. We said
 818 that a mutation is missense if the result of the genetic mutation is a different
 819 amino-acid in the protein. The mutation is nonsense when the genetic muta-
 820 tion results in a stop codon instead of an amino-acid. Finally, a silent muta-
 821 tion is a mutation with no effect on the amino-acid string, *i.e.* the mutated sequence
 822 codes for an amino-acid with identical binding properties. We talk about inser-
 823 tion (resp. deletion) when one or more nucleotides are added (resp. removed)
 824 at some place in the DNA code. These last kinds of mutations can both be
 825 frameshift mutations, which are given by the insertion or deletion of a number
 826 of bases that is not a multiple of 3, altering the reading frame of the gene.
 827 SHM introduces mostly single nucleotide exchanges, together with small dele-
 828 tions and duplications, *i.e.* the insertion of extra copies of a portion of genetic

829 material already present within the DNA code [35,14,15]. Among these point
830 mutations, transitions (*i.e.* substitution of a purine nucleotide with another
831 purine one, or a pyrimidine with a pyrimidine) dominate over transversions
832 (substitution of a purine with a pyrimidine or conversely). About half of the
833 mutations (53%) have been estimated to be silent, about 28% nonsense, and
834 only about 19% of all mutations have been estimated to be missense and then
835 have an effect on affinity, which can either be of an improving nature, or of
836 worsening and even lead to the formation of autoreactive clones [36].

837

838 The 20 existing amino-acids are typically classified in charged amino-acids,
839 polar (non-charged) amino-acids and hydrophobic amino-acids, depending on
840 their chemical characteristics. As we have already discussed in Section 4.2 the
841 bonding between BCR and antigen is made thanks to non-covalent bonds,
842 in particular ionic bonds and hydrogen bonds. Ionic bonds are the result of
843 interactions between two amino-acids oppositely charged: arginine (R) and
844 lysine (K) are positively charged, while aspartic acid (D) and glutamic acid
845 (E) are negatively charged. As long as hydrogen bonds are concerned, also
846 polar amino-acids can participate. In particular arginine (R), lysine (K) and
847 tryptophan (W) have hydrogen donor atoms in their side chains; aspartic acid
848 (D) and glutamic acid (E) have hydrogen acceptor atoms in their side chain
849 while asparagine (N), glutamine (Q), histidine (H), serine (S), threonine (T)
850 and tyrosine (Y) have both hydrogen donor and acceptor atoms in their side
851 chains.

852

853 Stop codons also have an important role. Indeed, during translation (the
854 last step necessary to build a protein starting from the DNA molecule) amino-
855 acids continue to be added until a stop codon is reached. There exists two
856 types of mutations involving stop codons, named nonsense and nonstop re-
857 spectively. The first one corresponds to the substitution of an amino-acid with
858 a stop codon, while the second one is the opposite case. In both cases the re-
859 sulting protein has an abnormal length, which often causes a loss of function.
860 Moreover, errors given by both nonsense and nonstop mutations are linked to
861 over 10% of human genetic diseases [12].

862

863 Concerning mutation in activated B-cells, SHM is driven by an enzyme
864 called activation-induced cytidine deaminase (AID) which is expressed specif-
865 ically in this case. This protein can bind to single-stranded DNA only. Thus
866 it seems to target only genes being transcribed (for which the transcription
867 phenomenon separates temporarily double stranded DNA into small portions
868 of two single stranded DNA sequences) [40]. AID converts Cytosine (C) in
869 Uracil (U) by deamination. This substitution occurs at higher rates in hot
870 spots motives like $D\underline{G}Y\underline{W}/\underline{W}R\underline{C}H$ where ($G : C$ is the mutable position and
871 $D \in \{A, G, T\}$, $H \in \{A, C, T\}$, $R \in \{A, G\}$, $W \in \{A, T\}$ and $Y \in \{C, T\}$, and the
872 underlined letters are the loci of mutations) [62,35]. Then, two mechanisms
873 tend to repair lesions in the DNA caused by these substitutions of C by U [63]:

- 874 a) either *mismatch repair* : substitution for the damaged zone by another
 875 sequence of nucleotides thanks to proteins MSH 2/6. The U base is read
 876 as T leading to a transition from a $C : G$ pair to $T : A$.
- 877 b) or *base excision repair* : U is excised by a successive action of uracil-
 878 DNA glycolase (UNG) and apurinic/aprimidinic endonuclease (APE1).
 879 The DNA contains then a nick, after replication, a random nucleotide is
 880 inserted in order to fill the vacant space leading to transversions and tran-
 881 sitions.

882 From a mathematical point of view this is equivalent to define the switch
 883 with a random nucleotide depending on the motives present in the chain. The
 884 probability concerning the choice of this nucleotide to be inserted shall not be
 885 uniform due to the presence of mismatch and excision repairs [20, 63]. This is
 886 not taken into account in the model we developed.

887
 888 We can therefore make the following three main assumptions to model the
 889 SHM process acting on the BCR V region:

890 *Modeling assumption 1* SHM introduces only single point mutations in the
 891 DNA strand, missense or silent. Therefore we do not take into account nonsense
 892 mutations, in order to avoid an interruption of the mutation process due to
 893 the introduction of a stop codon. The choice of the base used for substitution
 894 is made randomly, without considering that we have mostly $A \leftrightarrow T$ and $G \leftrightarrow C$
 895 substitutions.

896 *Modeling assumption 2* We consider only electrostatic and hydrogen bonds as
 897 responsible for the bonding between BCR and antigen. We suppose we have
 898 two amino-acid classes represented as 0 and 1 respectively: we denote by 1
 899 those amino-acids which have hydrogen donor atoms in their side chains (or
 900 which are positively charged) and by 0 those amino-acids which have hydrogen
 901 acceptor atoms in their side chains (or which are negatively charged). We
 902 arbitrary chose to assign 0 or 1 to amino-acids which can act as an acid or a
 903 base in hydrogen bonds. As an exemple, as serine can form hydrogen bonds
 904 with arginine and threonine, one can assign 0 to serine and 1 to threonine
 905 (arginine is represented by 1 as it is positively charged). While translating the
 906 amino-acid chain into a binary chain, we omit all hydrophobic amino-acids,
 907 as they do not participate in electrostatic or hydrogen bonds. Their position
 908 corresponds to an empty case, which does not contribute to the affinity between
 909 B-cell and antigen. This is clearly an important simplification. We will further
 910 discuss this choice in Section 4.4.

911 *Modeling assumption 3* We consider a linear contact between two amino-acid
 912 strings, without taking into account the geometrical configuration of both the
 913 BCR and the antigen.

914 The process starts from a DNA chain coding for a BCR, \mathbf{X}_0^{dna} ; from which
 915 we can obtain the corresponding amino-acid chain, \mathbf{X}_0^{aa} (Table 5) and, conse-
 916 quently, its binary expression, \mathbf{X}_0^{bin} .

917 *Example 1*

918
919 – $\mathbf{X}_0^{dna} = (\text{GTT, GAG, CTA, GTG, GAA, AGT, GGA, GCC, GAA, GTA, AAA,}$
920 $\text{AAG, CCA, GGT, AGT, AGT, GTT, AAA, GTC, AGT, TGT, AAA, GCA})$

921
922 – $\mathbf{X}_0^{aa} = (\text{V, Q, L, V, E, S, G, A, E, V, K, K, P, G, S, S, V, K, V, S, C, K, A})$

923
924 – $\mathbf{X}_0^{bin} = (-, 1, -, -, 0, 0, -, -, 0, -, 1, 1, -, -, 0, 0, -, 1, -, 0, 0, 1, -)$

925 *Notation 1* Given a vector \mathbf{X} , we denote by $|\mathbf{X}|$ its length (counting also the
926 empty cases, if there are some). Equivalently, given a set \mathcal{S} , we denote by $|\mathcal{S}|$
927 its size

928 We can formalize the translation of the nucleotides chain into the amino-
929 acids chain as follows.

930
931 **Definition 15** Let \mathcal{N} and \mathcal{A} be two sets of letters with size respectively $|\mathcal{N}| =$
932 k_1 and $|\mathcal{A}| = k_2$. Let l be an integer positive number so that $k_1^l \geq k_2$. Then
933 we define $f_{k_1, k_2, l} : \mathcal{N}^l \rightarrow \mathcal{A}$, which associate at least an l -length sequence of
934 letters belonging to \mathcal{N} to a letter in \mathcal{A} .

935 In our specific case, following definition 15, $\overline{\mathcal{N}} := \{\text{G, A, T, C}\}$ is the set
936 of nucleotides, while $\overline{\mathcal{A}}$ is the set containing all possible amino-acids, together
937 with the stop signal. Therefore $\overline{k}_1 = 4$ and $\overline{k}_2 = 21$. Moreover we know that
938 $\overline{l} = 3$ and the function $\overline{f}_{4, 21, 3}$ is detailed in Table 5.

939 *Remark 17* We can easily observe that $\overline{l} = \min \{n \in \mathbb{N} | \overline{k}_1^n \geq \overline{k}_2\}$. Indeed,
940 having 4 nucleotides available to build a DNA strand, we need to read them
941 at least by 3-length blocks in order to be able to synthesize all 20 amino-acids.
942 Moreover, choosing this value for the parameter l avoids to have too many
943 sequences of nucleotides coding for the same amino-acid.

944 At the beginning of the process, the antigen string in its three representa-
945 tions is given as well: $\overline{\mathbf{x}}^{dna}$, $\overline{\mathbf{x}}^{aa}$ and $\overline{\mathbf{x}}^{bin}$, with $|\mathbf{X}^{dna}| = |\overline{\mathbf{x}}^{dna}| =: 3N$. Anti-
946 gen representing strings remain unchanged. Assumptions 1-3 imply that for all
947 $t \geq 0$, $|\mathbf{X}_t^{bin}| = |\overline{\mathbf{x}}^{bin}| = N$. At each time step a single point mutation (missense
948 or silent) is introduced in the DNA chain coding for the BCR. So, if \mathbf{X}_t^{dna}
949 is the DNA code at time t , we randomly choose an index $i \in \{1, \dots, 3N\}$, a
950 letter $a \in \overline{\mathcal{N}}$ and we place $(X_{t+1}^{dna})_i := a$. If the new codon is a stop codon,
951 then we choose $a' \in \overline{\mathcal{N}} \setminus \{a\}$ and we put $(X_{t+1}^{dna})_i := a'$, and so on.
952 In order to test the affinity, we consider the binary expression of both the
953 BCR and the antigen, which we take in its complementary form, *i.e.* $\overline{\mathbf{x}}^{bin} :=$
954 $(1 - \overline{x}_1^{bin}, \dots, 1 - \overline{x}_N^{bin})$. This leads us back to the definition of affinity we made
955 in Section 2: 0 matches with 0 and 1 with 1.

956

957 As we consider a linear contact between \mathbf{X}_t^{bin} and $\bar{\mathbf{x}}'^{bin}$, at the positions
 958 where either \mathbf{X}_t^{bin} or $\bar{\mathbf{x}}'^{bin}$ has an hydrophobic amino-acid, we suppose that
 959 no match is possible. Therefore we can extend Definition 4 of the Hamming
 960 distance in a very natural way to this more general case:

Definition 16 We denote by $Hy(\mathbf{X}_t^{bin})$ (resp. $Hy(\bar{\mathbf{x}}'^{bin})$) the set of the indices corresponding to hydrophobic amino-acids in \mathbf{X}_t^{bin} (resp. in $\bar{\mathbf{x}}'^{bin}$). Therefore the Hamming distance between \mathbf{X}_t^{bin} and $\bar{\mathbf{x}}'^{bin}$ is given by:

$$h(\mathbf{X}_t^{bin}, \bar{\mathbf{x}}'^{bin}) = \sum_{\substack{i \in \{1, \dots, N\} \\ i \notin Hy(\mathbf{X}_t^{bin}) \cup Hy(\bar{\mathbf{x}}'^{bin})}} \delta_i + |Hy(\mathbf{X}_t^{bin}) \cup Hy(\bar{\mathbf{x}}'^{bin})|$$

$$\text{where } \delta_i = \begin{cases} 1 & \text{if } (X_t^{bin})_i \neq (\bar{x}'^{bin})_i \\ 0 & \text{otherwise} \end{cases}$$

Then, for all $t \geq 0$:

$$|Hy(\mathbf{X}_t^{bin}) \cup Hy(\bar{\mathbf{x}}'^{bin})| \leq h(\mathbf{X}_t^{bin}, \bar{\mathbf{x}}'^{bin}) \leq N$$

We consider that the optimal clone is reached when:

$$\text{aff}(\mathbf{X}_t^{bin}, \bar{\mathbf{x}}'^{bin}) := N - |Hy(\bar{\mathbf{x}}'^{bin})|$$

961 The effects of nucleotides exchanges on the binary expression of BCRs can
 962 be multiple:

963 **No detectable effect** : this is the result of either a silent mutation or a mis-
 964 sense mutation which substitutes an amino-acid with another one belonging
 965 to the same amino-acid class.

966 **Class-switch** , derived from a missense mutation which leads to the substitu-
 967 tion of an amino-acid with another one belonging to the other amino-acid
 968 class.

969 We can further complexify this model by replacing Assumption 1 with the
 970 following one:

971 *Modeling assumption 4* SHM introduces mostly single point mutations in the
 972 DNA, missense or silent. With weak probability, deletions or insertions can
 973 occur. For the sake of simplicity, we suppose that a deletion (resp. an insertion)
 974 consist in the elimination (resp. the addition) of a non-stop codon. Moreover, in
 975 order to avoid the problem of a variation in the length of the BCR representing
 976 string, when a deletion occur, those bits situated on the right of the deleted
 977 one shift to the left, and a random extra codon is added at the right bottom.
 978 Conversely, if an insertion occurs, the right bottom bit is deleted.

979 Even if these mutational events are rare, they have remarkable effects over
 980 the structure of the underlying graph. Indeed a deletion or an insertion entails
 981 a great jump in the affinity function by producing a shift of a portion of the
 982 BCR representing string. This is not the case if we consider only single point
 983 mutations. Therefore, under Assumption 4 the graph we obtain is much more
 984 complex and allows random long range connections.

985 *4.3.1 Numerical simulations*

986 In order to evaluate how deletions and insertions affect the mean number of
 987 mutation steps to reach the desired B-cell trait, we make some numerical simu-
 988 lations. We compare a model in which only single point mutations are allowed
 989 to another one in which also deletions and insertions can occur. We refer to
 990 Assumption 4 to define these mutational events.

991

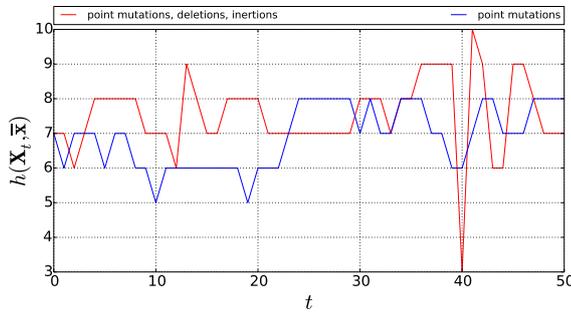


Figure 4: Variation of the Hamming distance to \bar{x}^{bin} , comparing the model of single point mutations to the one which includes also deletions and insertions (50% of all mutation events). In both cases $N = 10$. Deletions and insertions lead to a quick change in the Hamming distance. Between time 30 and 50, we can observe the effect of indels mutations.

992 Figure 4 shows the effects of deletions and insertions over the affinity. In
 993 order to do these simulations, we arbitrary fixe a BCR and an antigen with
 994 given affinity. We do not consider those base substitutions leading to no de-
 995 tectable effect, *i.e.* at each time step we can observe a variation of the affinity
 996 function. In Figure 4 we can clearly locate at what time an insertion or a
 997 deletion has occurred, because this coincides with a jump of the Hamming
 998 distance between BCR and antigen.

999

1000 One can ask how these random long range connections affect the average
 1001 time to reach the antigen target string. Simulations show that one needs a
 1002 more long time to reach \bar{x}^{bin} if the probability of making such mutations in-
 1003 creases. The results obtained through 10000 simulations are collected in Table

1004 **6.**
1005

Table 6: Average number of mutations needed to reach \bar{x}^{bin} , for $N = 10$ and starting from Hamming distance 7. In \bar{x}^{bin} , only 2 amino-acids are hydrophobic, so by Definition 16, the optimal affinity one can reach is 8. We compare three models: in the first one no deletions nor insertions are allowed. In the second model 10% of all mutations are deletions or insertions, 50% in the last one. We denote by $\widehat{\tau}_{\{\bar{x}^{bin}\}_n}$ the average value obtained over n simulations and by $\widehat{\sigma}_n$ its corresponding estimated standard deviation. Simulations show that $\widehat{\tau}_{\{\bar{x}^{bin}\}_n}$ increases when the pourcentage of deletions or insertions grows, and so does the corresponding variation.

% deletions/insertions	$ \bar{x}^{bin} $	$h(\mathbf{X}_0^{bin}, \bar{x}^{bin})$	n	$\widehat{\tau}_{\{\bar{x}^{bin}\}_n}$	$\widehat{\frac{\sigma_n}{\sqrt{n}}}$
0	10	7	10000	8824.93	86.80
10	10	7	10000	9091.12	92.01
50	10	7	10000	10075.89	100.59

1006 We can discuss which viewpoint is the most suitable to study mutations
1007 and their effects over the interactions between BCR and antigen. It is really
1008 hard to define a clear correspondence between genetic mutations and the evo-
1009 lution of the affinity, even while considering a simple linear contact between
1010 molecules (hence without observing the changes in the geometrical structure
1011 of the protein). Indeed, in order to test the affinity between BCR and anti-
1012 gen we constantly need to project the DNA string on the smaller state-space
1013 containing the binary representations of B-cell traits. If we directly consider
1014 mutations on binary strings, then the resulting process is faster, as we do not
1015 observe missense mutations, and the evaluation of the affinity is immediate.
1016

1017 The comprehension of the nature of genetic mutations and their conse-
1018 quences on the new generated protein, suggested us to make Assumptions 1-3
1019 to formalize the model. In particular, we found reasonable to look directly
1020 to amino-acid chains and their binary representation: this allows to study the
1021 affinity between BCR and antigen using the Hamming distance. Therefore, un-
1022 der these hypotheses the general mathematical framework described in Section
1023 2 can be applied to study how different kinds of missense mutations affect the
1024 dynamics of AAM. As we show in Sections 2-3, this already brings interesting
1025 and complexes mathematical problems.

1026 4.4 Limitations and extensions

1027 In this paper we propose and study mutational processes on N -length binary
1028 strings, which can be variously applied to evolutionary contexts. As far as
1029 the application to the SHM process is concerned, we can make some remarks
1030 about our assumptions, which can bring us to enrich and complexify the model
1031 through a more coherent representation of the true biological process.

1032
1033 First of all we have decided to consider only two amino-acid classes. From
1034 one side this assumption is justified as charged and polar amino-acids are ef-
1035 fectively the most responsible in creating bonds which determine the antigen-
1036 antibody interaction. Therefore they strongly influence the affinity between
1037 BCR and antigen. Nevertheless, by making this simplification we omit all
1038 hydrophobic amino-acids from the string, and that is not without conse-
1039 quences. The elimination of hydrophobic amino-acids from the string signif-
1040 icantly changes the structure of the chain, therefore the ability for charged
1041 and polar amino-acids to be in contact with each-others. Moreover, the effects
1042 of genetic mutations on the new generated protein could be even more com-
1043 plex than the ones we have considered in this paper. Finally, by taking into
1044 account also hydrophobic amino-acids, we would be able to consider hydropho-
1045 bic bonds, which also influences the antigen-antibody interaction. Therefore it
1046 seems more appropriate to consider three, or more, amino-acids classes (*e.g.*
1047 [\[59, 53\]](#)).

1048
1049 As far as the nature of mutations is concerned, we have essentially de-
1050 scribed mutational processes given by combinations of single point mutation
1051 mechanisms. During SHM nucleotide exchanges are the most frequent among
1052 all possible mutations. Despite this, also some deletions and insertions occur.
1053 This has two main consequences. Firstly it means that the length of the BCR
1054 representing string could change during the process, while we consider it as
1055 fixed and equal to the length of the antigen. We can maybe overcome this
1056 problem by saying that the chain represented in our model corresponds to the
1057 portion of BCR in contact with the antigen, and this is almost fixed (Section
1058 [4.2](#)). Moreover these mutations can imply substantial changes into the amino-
1059 acid chain, hence they can bring a great jump of the affinity to the presented
1060 antigen. Therefore, even if these are rare mutational events, they may have
1061 an important effect in AAM. Consequently it could be interesting to take also
1062 insertions and deletions into account. All these observations lead interesting
1063 mathematical questions.

1064
1065 Of course we can also envisage developments in other directions. For exam-
1066 ple by considering the creation of bonds among amino-acids of the BCR (resp.
1067 the antigen) itself, which determines the geometrical structure of the protein
1068 and consequently the portion of the BCR and the antigen that can actually be
1069 in contact. Another interesting possibility is to consider that mutations at one
1070 site are influenced by other amino acids composing the string. This assumption

1071 was firstly proposed by S. A. Kauffman and E. D. Weinberger in [39], where
 1072 they introduced the NK models. In this context the parameter K assures the
 1073 richness of epistatic interactions among sites. More recently Y. Elhanati *et al*
 1074 in [23] find biological evidence for an evolutionary model where substitution
 1075 rates strictly depend on the context.

1076
 1077 We propose some numerical simulations to evaluate the consequences over
 1078 the hitting time of both the addition of extra amino-acid classes and the pos-
 1079 sibility of having a BCR string longer than the antigen one.

1080
 1081 A. S. Perelson and G. Weisbuch in [59] proposed a model with 3 amino-
 1082 acid classes: hydrophobic, hydrophilic positively charged and hydrophilic neg-
 1083 atively charged. Hydrophobic amino-acids match with hydrophobic and hy-
 1084 drophilic positively charged with hydrophilic negatively charged. We simulated
 1085 the expected time to reach a given configuration comparing the model with
 1086 2 amino-acid classes and the one with 3 amino-acid classes, and considering
 1087 single switch-type mutations. We take two random 10-length strings having
 1088 maximal distance between each-others. We extend Definition 4 of Hamming
 1089 distance to the state-space $\{0, 1, 2\}^N$ in a natural way, keeping the same nota-
 1090 tion: $\forall \mathbf{x} = (x_1, \dots, x_N), \mathbf{y} = (y_1, \dots, y_N) \in \{0, 1, 2\}^N$, their Hamming distance
 1091 is given by:

$$h(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^N \delta_i \quad \text{where} \quad \delta_i = \begin{cases} 1 & \text{if } x_i \neq y_i \\ 0 & \text{otherwise} \end{cases} \quad (31)$$

1092 Therefore the affinity is defined as in Definition 3. We simulated for both cases
 1093 a single switch-type mutational model (Definition 5 for 2 amino-acid classes
 1094 and Definition 17 below for 3 amino-acid classes), testing the time we need to
 1095 reach the target vertex.

1096 **Definition 17** Let $\mathbf{X}_n \in \{0, 1, 2\}^N$ be the BCR at step n . Let $i \in \{1, \dots, N\}$
 1097 be a randomly chosen index, and $a \in \{0, 1, 2\} \setminus \{X_{n,i}\}$ a randomly chosen
 1098 number. Then $\mathbf{X}_{n+1} := (X_{n,1}, \dots, X_{n,i-1}, a, X_{n,i+1}, \dots, X_{n,N})$.

1099 Table 7 shows the results we obtained over 10000 simulations.

1100 We already knew from theoretical analysis that the order of magnitude
 1101 for the hitting time of the basic mutational model is 2^N for N big enough.
 1102 Simulations clearly show that when we consider 3 amino-acid classes, the or-
 1103 der of magnitude of the hitting time of a single switch-type mutational model
 1104 significantly increases, and is of the order of 3^N , as proved by Proposition 4.
 1105 Moreover we observe that the variance corresponding to the second model is
 1106 significantly bigger as well.

1107
 1108 It is clear that if we consider more amino-acid classes, it takes much longer
 1109 to reach a precise element of the new state-space. Nevertheless, one can under-
 1110 stand that if we keep the same distance function as defined in Equation (31),

Table 7: Average expected times to cover a Hamming distance $h(\mathbf{X}_0, \bar{\mathbf{x}}) = 10 = N$, comparing the model with 2 amino-acid classes and the one with 3 amino-acid classes. Here we denote by $\widehat{\tau_{\{\bar{\mathbf{x}}\}}_n}$ the average value obtained over n simulations and by $\widehat{\sigma}_n$ its corresponding estimated standard deviation.

Amino-acid classes	N	$h(\mathbf{X}_0, \bar{\mathbf{x}})$	n	$\widehat{\tau_{\{\bar{\mathbf{x}}\}}_n}$	$\widehat{\frac{\sigma_n}{\sqrt{n}}}$
2	10	10	10000	1213.2108	12.0138
3	10	10	10000	62160.8263	635.0458

1111 than we are asking for a higher degree of precision while building the B-cell
 1112 trait. Therefore, we can not directly compare hitting times corresponding to
 1113 a model with a greater number of amino-acid classes and keeping the same
 1114 affinity function as the one used with only two amino-acid classes. If one want
 1115 to obtain a comparable result by using more than two amino-acid classes, one
 1116 has to use a weaker definition of affinity.

Definition 18 Let \mathcal{S} be a set of letters, $|\mathcal{S}| = s > 2$. Let us partition \mathcal{S} into two subsets: $\mathcal{S} := \mathcal{S}_1 \sqcup \mathcal{S}_2$. $\forall \mathbf{x}, \mathbf{y} \in \mathcal{S}^N$, their distance is given by:

$$h_{\mathcal{S}_1, \mathcal{S}_2}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^N \delta_i \quad \text{where} \quad \delta_i = \begin{cases} 1 & \text{if } x_i \in \mathcal{S}_1, y_i \in \mathcal{S}_2 \text{ or conversely} \\ 0 & \text{otherwise} \end{cases}$$

Consequently, their affinity is given by:

$$\text{aff}(\mathbf{x}, \mathbf{y}) = N - h_{\mathcal{S}_1, \mathcal{S}_2}(\mathbf{x}, \mathbf{y})$$

1117 By using this new affinity function we can compare the hitting times and
 1118 the order of magnitude is clearly the same.

1119

1120 Let us now go back to Assumption 2 and to the structure of the string
 1121 given in Section 4.3 (in particular, hydrophobic amino-acids are represented
 1122 by empty cases). Contrary to what stated by Assumption 4, we suppose that
 1123 the BCR length can be modified by insertions and deletions. Consequently, also
 1124 a modification of the distance function is needed. We arbitrarily fixe a BCR
 1125 and an antigen with given affinity. We do not consider those base substitutions
 1126 leading to no detectable effect, *i.e.* at each time step we can observe a variation
 1127 of the affinity function. We suppose that 90% of all mutation events are single
 1128 point mutations, 10% deletions or insertions. If we are in this case and $|\mathbf{X}_t^{bin}| >$
 1129 $|\bar{\mathbf{x}}^{bin}|$, then with probability 1/2 a deletion occurs and with probability 1/2 an
 1130 insertion occur. Otherwise, it will be necessarily an insertion (this is to avoid
 1131 to obtain $|\mathbf{X}_t^{bin}| = 0$). As long as the affinity is concerned, if $|\mathbf{X}_t^{bin}| > |\bar{\mathbf{x}}^{bin}|$,

1132 $|\mathbf{X}_t^{bin}| := n_1$, $|\bar{\mathbf{x}}'^{bin}| := n_2$, then their distance is the smaller possible one, *i.e.*:

$$h(\mathbf{X}_t^{bin}, \bar{\mathbf{x}}'^{bin}) = \min_{1 \leq i \leq n_1 - n_2 + 1} \left\{ h(\mathbf{X}_i, \bar{\mathbf{x}}'^{bin}) \mid \mathbf{X}_i := \left(X_{t,i}^{bin}, X_{t,i+1}^{bin}, \dots, X_{t,i+n_2-1}^{bin} \right) \right\},$$

1133 h as in Definition 16.

Table 8: Average number of mutations needed to reach $\bar{\mathbf{x}}'^{bin}$, for $N = 7$ and starting from a Hamming distance 5. In $\bar{\mathbf{x}}'^{bin}$, only 2 amino-acids are hydrophobic, so by Definition 16, the optimal Hamming distance one can reach is 2. We compare a model in which no deletions nor insertions are allowed and a model in which 10% of all mutations are deletions or insertions. We denote by $\tau_{\{\bar{\mathbf{x}}'^{bin}\}_n}$ the average value obtained over n simulations and by $\hat{\sigma}_n$ its corresponding estimated standard deviation.

% deletions/insertions	$ \bar{\mathbf{x}}'^{bin} $	$h(\mathbf{X}_0^{bin}, \bar{\mathbf{x}}'^{bin})$	n	$\tau_{\{\bar{\mathbf{x}}'^{bin}\}_n}$	$\frac{\hat{\sigma}_n}{\sqrt{n}}$
0	7	5	5000	374.28	5.38
10	7	5	5000	251.48	3.54

1133 In this case, and thanks to the definition of Hamming distance as the min-
 1134 imal one, we clearly have more chances to obtain a good B-cell trait. This is
 1135 confirmed by results collected in Table 8. When deletions and insertions can
 1136 occur, even with very weak probability, and if we allowed the BCR length
 1137 to be greater than the antigen one, then the expected number of mutations
 1138 needed to built the optimal BCR is more than 30% smaller.

1139

1140 5 Conclusion

1141 In this paper, we have introduced a mathematical framework to study the
 1142 impact of various mutation rules on the exploration of the space of traits in an
 1143 evolutionary model. In particular, we have connected mutation rules to char-
 1144 acteristic time-scales, such as hitting-times, through the study of associated
 1145 graph structures. As a leading example, which was the original motivation for
 1146 this study, we have considered applications of these results to the modeling of
 1147 somatic hypermutations in the germinal center. The models considered so far
 1148 do not include division and selection, which would lead to studying branching
 1149 random walks on graphs, a topic of ongoing research.

1150 References

- 1151 1. Abbas, A.K., Lichtman, A.H., Pillai, S.: Basic immunology: functions and disorders of
 1152 the immune system. Elsevier Health Sciences (2012)

- 1153 2. Aickelin, U., Dasgupta, D., Gu, F.: Artificial immune systems. In: Search Methodologies,
1154 pp. 187–211. Springer (2014)
- 1155 3. Aldous, D., Fill, J.: Reversible markov chains and random walks on graphs (2002)
- 1156 4. Ansari, H.R., Raghava, G.P.: Identification of conformational b-cell epitopes in an anti-
1157 gen from its primary sequence. *Immunome research* **6**(1), 1 (2010)
- 1158 5. Bäck, T.: Evolutionary algorithms in theory and practice: evolution strategies, evolu-
1159 tionary programming, genetic algorithms. Oxford university press (1996)
- 1160 6. Balelli, I., Milisic, V., Wainrib, G.: Branching random walks on binary strings for evo-
1161 lutionary processes. arXiv preprint arXiv:1607.00927 (2016)
- 1162 7. Balelli, I., Milišić, V., Wainrib, G.: Multi-type galton-watson processes with affinity-
1163 dependent selection applied to antibody affinity maturation. arXiv preprint
1164 arXiv:1609.00823 (2016)
- 1165 8. Benson, M.J., Erickson, L.D., Gleeson, M.W., Noelle, R.J.: Affinity of antigen encounter
1166 and other early b-cell signals determine b-cell fate. *Current opinion in immunology*
1167 **19**(3), 275–280 (2007)
- 1168 9. Berestycki, N.E.: Phase transitions for the distance of random walks with applications
1169 to genome rearrangements. Ph.D. thesis, Cornell University (2005)
- 1170 10. Berkhin, P.: A survey on pagerank computing. *Internet Mathematics* **2**(1), 73–120
1171 (2005)
- 1172 11. Besmer, E., Gourzi, P., Papavasiliou, F.N.: The regulation of somatic hypermutation.
1173 *Current opinion in immunology* **16**(2), 241–245 (2004)
- 1174 12. Bidou, L., Allamand, V., Rousset, J.P., Namy, O.: Sense from nonsense: therapies for
1175 premature stop codon diseases. *Trends in molecular medicine* **18**(11), 679–688 (2012)
- 1176 13. Binitha, S., Sathya, S.S.: A survey of bio inspired optimization algorithms. *International*
1177 *Journal of Soft Computing and Engineering* **2**(2), 137–151 (2012)
- 1178 14. Bowers, P.M., Verdino, P., Wang, Z., da Silva Correia, J., Chhoa, M., Macondray, G.,
1179 Do, M., Neben, T.Y., Horlick, R.A., Stanfield, R.L., et al.: Nucleotide insertions and
1180 deletions complement point mutations to massively expand the diversity created by
1181 somatic hypermutation of antibodies. *Journal of Biological Chemistry* **289**(48), 33,557–
1182 33,567 (2014)
- 1183 15. Briney, B.S., Willis, J.R., Crowe, J.: Location and length distribution of somatic
1184 hypermutation-associated dna insertions and deletions reveals regions of antibody struc-
1185 tural plasticity. *Genes and immunity* **13**(7), 523–529 (2012)
- 1186 16. Castro, L.N.D., Zuben, F.J.V.: Learning and optimization using the clonal selection
1187 principle. *Evolutionary Computation, IEEE Transactions on* **6**(3), 239–251 (2002)
- 1188 17. Cobey, S., Wilson, P., Matsen, F.A.: The evolution within us. *Phil. Trans. R. Soc. B*
1189 **370**(1676), 20140,235 (2015)
- 1190 18. Currin, A., Swainston, N., Day, P.J., Kell, D.B.: Synthetic biology for the directed
1191 evolution of protein biocatalysts: navigating sequence space intelligently. *Chemical*
1192 *Society Reviews* **44**(5), 1172–1239 (2015)
- 1193 19. De Silva, N.S., Klein, U.: Dynamics of b cells in germinal centres. *Nature Reviews*
1194 *Immunology* **15**(3), 137–148 (2015)
- 1195 20. Di Noia, J., Neuberger, M.S.: Altering the pathway of immunoglobulin hypermutation
1196 by inhibiting uracil-DNA glycosylase. *Nature* **419**(6902), 43–48 (2002)
- 1197 21. Diaconis, P., Graham, R.L., Morrison, J.A.: Asymptotic analysis of a random walk on a
1198 hypercube with many dimensions. *Random Structures & Algorithms* **1**(1), 51–72 (1990)
- 1199 22. Doyle, P.G., Snell, J.L.: Random walks and electric networks. *AMC* **10**, 12 (1984)
- 1200 23. Elhanati, Y., Sethna, Z., Marcou, Q., Callan, C.G., Mora, T., Walczak, A.M.: Infer-
1201 ring processes underlying b-cell repertoire diversity. *Phil. Trans. R. Soc. B* **370**(1676),
1202 20140,243 (2015)
- 1203 24. Ethier, S.N., Kurtz, T.G.: Markov processes: characterization and convergence, vol. 282.
1204 John Wiley & Sons (2009)
- 1205 25. Ewens, W.J.: Mathematical population genetics. i. theoretical introduction. *interdisci-
1206 plinary applied mathematics*, 27 (2004)
- 1207 26. Faro, J., Or-Guil, M.: How oligoclonal are germinal centers? a new method for estimating
1208 clonal diversity from immunohistological sections. *BMC bioinformatics* **14**(Suppl 6), S8
1209 (2013)
- 1210 27. Fisher, R.A.: The genetical theory of natural selection: a complete variorum edition.
1211 Oxford University Press (1930)

- 1212 28. Florkowski, S.F.: Spectral graph theory of the hypercube. Master's thesis, Naval Post-
1213 graduate School, Monterey, California (2008)
- 1214 29. Forrest, R.E.S.S., Perelson, A.S.: Population diversity in an immune system model:
1215 Implications for genetic search. *Foundations of Genetic Algorithms 1993 (FOGA 2)* **2**,
1216 153 (2014)
- 1217 30. Frost, S.D., Murrell, B., Hossain, A.M.M., Silverman, G.J., Pond, S.L.K.: Assigning and
1218 visualizing germline genes in antibody repertoires. *Phil. Trans. R. Soc. B* **370**(1676),
1219 20140,240 (2015)
- 1220 31. Gitlin, A.D., Shulman, Z., Nussenzweig, M.C.: Clonal selection in the germinal centre
1221 by regulated proliferation and hypermutation. *Nature* (2014)
- 1222 32. Gjoka, M., Kurant, M., Butts, C.T., Markopoulou, A.: Walking in facebook: A case
1223 study of unbiased sampling of osns. In: *INFOCOM, 2010 Proceedings IEEE*, pp. 1–9.
1224 *IEEE* (2010)
- 1225 33. Haldane, J.B.S.: The cost of natural selection. *Journal of Genetics* **55**(3), 511–524 (1957)
- 1226 34. Harary, F., Hayes, J.P., Wu, H.J.: A survey of the theory of hypercube graphs. *Com-
1227 puters & Mathematics with Applications* **15**(4), 277–289 (1988)
- 1228 35. Hwang, J.K., Alt, F.W., Yeap, L.S.: Related mechanisms of antibody somatic hyper-
1229 mutation and class switch recombination. *Microbiology spectrum* **3**(1) (2015)
- 1230 36. Iber, D., Maini, P.K.: A mathematical model for germinal centre kinetics and affinity
1231 maturation. *Journal of theoretical biology* **219**(2), 153–175 (2002)
- 1232 37. Jeh, G., Widom, J.: Simrank: a measure of structural-context similarity. In: *Proceedings
1233 of the eighth ACM SIGKDD international conference on Knowledge discovery and data
1234 mining*, pp. 538–543. *ACM* (2002)
- 1235 38. Kamp, C., Bornholdt, S.: Coevolution of quasispecies: B-cell mutation rates maximize
1236 viral error catastrophes. *Physical Review Letters* **88**(6), 068,104 (2002)
- 1237 39. Kauffman, S.A., Weinberger, E.D.: The nk model of rugged fitness landscapes and its
1238 application to maturation of the immune response. *Journal of theoretical biology* **141**(2),
1239 211–245 (1989)
- 1240 40. Keim, C., Kazadi, D., Rothschild, G., Basu, U.: Regulation of AID, the B-cell genome
1241 mutator. *Genes Dev.* **27**(1), 1–17 (2013)
- 1242 41. Kempe, D., Dobra, A., Gehrke, J.: Gossip-based computation of aggregate information.
1243 In: *Foundations of Computer Science, 2003. Proceedings. 44th Annual IEEE Symposium
1244 on*, pp. 482–491. *IEEE* (2003)
- 1245 42. Kepler, T.B., Perelson, A.S.: Cyclic re-entry of germinal center b cells and the efficiency
1246 of affinity maturation. *Immunology today* **14**(8), 412–415 (1993)
- 1247 43. Kepler, T.B., Perelson, A.S.: Somatic hypermutation in b cells: an optimal control treat-
1248 ment. *Journal of theoretical biology* **164**(1), 37–64 (1993)
- 1249 44. Konstas, I., Stathopoulos, V., Jose, J.M.: On social networks and collaborative recom-
1250 mendation. In: *Proceedings of the 32nd international ACM SIGIR conference on
1251 Research and development in information retrieval*, pp. 195–202. *ACM* (2009)
- 1252 45. Kringelum, J.V., Lundegaard, C., Lund, O., Nielsen, M.: Reliable b cell epitope pre-
1253 dictions: impacts of method development and improved benchmarking. *PLoS Comput
1254 Biol* **8**(12), e1002,829 (2012)
- 1255 46. Kringelum, J.V., Nielsen, M., Padkjær, S.B., Lund, O.: Structural analysis of b-cell
1256 epitopes in antibody: protein complexes. *Molecular immunology* **53**(1), 24–34 (2013)
- 1257 47. Krovi, H., Brun, T.A.: Hitting time for quantum walks on the hypercube. *Physical
1258 Review A* **73**(3), 032,341 (2006)
- 1259 48. Levin, D.A., Peres, Y., Wilmer, E.L.: *Markov chains and mixing times*. Amer Mathe-
1260 matical Society (2009)
- 1261 49. Lovász, L.: Random walks on graphs: A survey. *Combinatorics, Paul erdos is eighty*
1262 **2**(1), 1–46 (1993)
- 1263 50. Meyer-Hermann, M.: A mathematical model for the germinal center morphology and
1264 affinity maturation. *Journal of theoretical Biology* **216**(3), 273–300 (2002)
- 1265 51. Meyer-Hermann, M., Mohr, E., Pelletier, N., Zhang, Y., Vitorica, G.D., Toellner, K.M.:
1266 A theory of germinal center b cell selection, division, and exit. *Cell reports* **2**(1), 162–174
1267 (2012)
- 1268 52. Meyn, S.P., Tweedie, R.L.: *Markov chains and stochastic stability*. Cambridge University
1269 Press (2009)

- 1270 53. Muñoz, E., Deem, M.W.: Amino acid alphabet size in protein evolution experiments:
1271 better to search a small library thoroughly or a large library sparsely? *Protein Engi-*
1272 *neering Design and Selection* **21**(5), 311–317 (2008)
- 1273 54. Murphy, K.M., Travers, P., Walport, M., et al.: *Janeway’s immunobiology*, vol. 7. Gar-
1274 *land Science* New York, NY, USA (2012)
- 1275 55. Norris, J.R.: *Markov chains*. 2008. Cambridge university press (1998)
- 1276 56. Oprea, M., Perelson, A.S.: Somatic mutation leads to efficient affinity maturation when
1277 centrocytes recycle back to centroblasts. *The Journal of Immunology* **158**(11), 5155–
1278 5162 (1997)
- 1279 57. Or-Guil, M., Faro, J.: A major hindrance in antibody affinity maturation investigation:
1280 We never succeeded in falsifying the hypothesis of single-step selection. *Frontiers in*
1281 *Immunology* **5** (2014)
- 1282 58. Pang, W., Wang, K., Wang, Y., Ou, G., Li, H., Huang, L.: Clonal selection algorithm for
1283 solving permutation optimisation problems: A case study of travelling salesman prob-
1284 lem. In: *International Conference on Logistics Engineering, Management and Computer*
1285 *Science (LEMCS 2015)*. Atlantis Press (2015)
- 1286 59. Perelson, A.S., Weisbuch, G.: *Immunology for physicists*. *Reviews of modern physics*
1287 **69**(4), 1219–1267 (1997)
- 1288 60. Ramiscal, R.R., Vinuesa, C.G.: T-cell subsets in the germinal center. *Immunological*
1289 *reviews* **252**(1), 146–155 (2013)
- 1290 61. Rogers, L.C.G., Williams, D.: *Diffusions, Markov Processes, and Martingales: Volume*
1291 *1, Foundations*. Cambridge university press (2000)
- 1292 62. Rogozin, I.B., Diaz, M.: Cutting edge: DGYW/WRCH is a better predictor of mutability
1293 at G:C bases in Ig hypermutation than the widely accepted RGYW/WRCY motif and
1294 probably reflects a two-step activation-induced cytidine deaminase-triggered process. *J.*
1295 *Immunol.* **172**(6), 3382–3384 (2004)
- 1296 63. Saribasak, H., Gearhart, P.J.: Does dna repair occur during somatic hypermutation?
1297 In: *Seminars in immunology*, 4, pp. 287–292. Elsevier (2012)
- 1298 64. Schaeffer, S.E.: Graph clustering. *Computer Science Review* **1**(1), 27–64 (2007)
- 1299 65. Sciammas, R., Li, Y., Warmflash, A., Song, Y., Dinner, A.R., Singh, H.: An incoherent
1300 regulatory network architecture that orchestrates b cell diversification in response to
1301 antigen signaling. *Molecular systems biology* **7**(1), 495 (2011)
- 1302 66. Shannon, M., Mehr, R.: Reconciling repertoire shift with affinity maturation: the role
1303 of deleterious mutations. *The Journal of Immunology* **162**(7), 3950–3956 (1999)
- 1304 67. Shen, W.J., Wong, H.S., Xiao, Q.W., Guo, X., Smale, S.: Towards a mathematical
1305 foundation of immunology and amino acid chains. *arXiv preprint arXiv:1205.6031* (2012)
- 1306 68. Shulman, Z., Gitlin, A.D., Targ, S., Jankovic, M., Pasqual, G., Nussenzweig, M.C.,
1307 Victora, G.D.: T follicular helper cell dynamics in germinal centers. *Science* **341**(6146),
1308 673–677 (2013)
- 1309 69. Smith, A.: Nucleic acids to amino acids: Dna specifies protein. *Nature Education* **1**(1),
1310 126 (2008)
- 1311 70. Sompayrac, L.: *How the immune system works*. Wiley-Blackwell (2012)
- 1312 71. Stern, J.N., O’Connor, K.C., Hafler, D.A., Laserson, U., Vigneault, F., Kleinstein, S.H.:
1313 Models of somatic hypermutation targeting and substitution based on synonymous muta-
1314 tions from high-throughput immunoglobulin sequencing data. *Immune system mod-*
1315 *eling and analysis* p. 55 (2015)
- 1316 72. Tas, J.M., Mesin, L., Pasqual, G., Targ, S., Jacobsen, J.T., Mano, Y.M., Chen, C.S.,
1317 Weill, J.C., Reynaud, C.A., Browne, E.P., Meyer-Hermann, M., Victora, G.D.: Visual-
1318 izing antibody affinity maturation in germinal centers. *Science* **351**(6277), 1048–1054
1319 (2016)
- 1320 73. Teng, G., Papavasiliou, F.N.: Immunoglobulin somatic hypermutation. *Annu. Rev.*
1321 *Genet.* **41**, 107–120 (2007)
- 1322 74. Tonegawa, S.: Somatic generation of immune diversity. *Bioscience reports* **8**(1), 3–26
1323 (1988)
- 1324 75. Victora, G.D.: Snapshot: the germinal center reaction. *Cell* **159**(3), 700–700 (2014)
- 1325 76. Victora, G.D., Schwickert, T.A., Fooksman, D.R., Kamphorst, A.O., Meyer-Hermann,
1326 M., Dustin, M.L., Nussenzweig, M.C.: Germinal center dynamics revealed by multi-
1327 photon microscopy with a photoactivatable fluorescent reporter. *Cell* **143**(4), 592–605
1328 (2010)

-
- 1329 77. Voit, M.: Asymptotic distributions for the ehrenfest urn and related random walks.
1330 Journal of Applied Probability pp. 340–356 (1996)
- 1331 78. Wang, F., Sen, S., Zhang, Y., Ahmad, I., Zhu, X., Wilson, I.A., Smider, V.V., Magliery,
1332 T.J., Schultz, P.G.: Somatic hypermutation maintains antibody thermodynamic sta-
1333 bility during affinity maturation. Proceedings of the National Academy of Sciences
1334 **110**(11), 4261–4266 (2013)
- 1335 79. Wright, S.: The roles of mutation, inbreeding, crossbreeding, and selection in evolution,
1336 vol. 1. na (1932)