

1 A mathematical model clarifies the ABC Score formula used in 2 enhancer-gene prediction

3 Joseph Nasser^{1,2,*}, Kee-Myoung Nam^{1,3}, Jeremy Gunawardena^{1,*}

4 ¹Department of Systems Biology, Harvard Medical School, Boston, MA, USA

5 ²Current address: Department of Physics, Brandeis University, Waltham, MA, USA

6 ³Current address: Department of Molecular, Cellular and Developmental Biology, Yale University, New
7 Haven, CT, USA

8 *To whom correspondence may be addressed: Joseph Nasser (joseph.nasser21@gmail.com) and Jeremy
9 Gunawardena (jeremy@hms.harvard.edu)

10 Abstract

11 Enhancers are discrete DNA elements that regulate the expression of eukaryotic genes. They are important
12 not only for their regulatory function, but also as loci that are frequently associated with disease traits.
13 Despite their significance, our conceptual understanding of how enhancers work remains limited. CRISPR-
14 interference methods have recently provided the means to systematically screen for enhancers in cell culture,
15 from which a formula for predicting whether an enhancer regulates a gene, the Activity-by-Contact (ABC)
16 Score, has emerged and has been widely adopted. While useful as a binary classifier, it is less effective at
17 predicting the quantitative effect of an enhancer on gene expression. It is also unclear how the algebraic
18 form of the ABC Score arises from the underlying molecular mechanisms and what assumptions are needed
19 for it to hold. Here, we use the graph-theoretic linear framework, previously introduced to analyze gene
20 regulation, to formulate the *default model*, a mathematical model of how multiple enhancers independently
21 regulate a gene. We show that the algebraic form of the ABC Score arises from this model. However, the
22 default model assumptions also imply that enhancers act additively on steady-state gene expression. This
23 is known to be false for certain genes and we show how modifying the assumptions can accommodate this
24 discrepancy. Overall, our approach lays a rigorous, biophysical foundation for future studies of enhancer-gene
25 regulation.

26 Introduction

27 Much of our current understanding of how genes are regulated arose from classical studies in bacteria of the
28 lac operon and λ -phage [1]. However, the eukaryotic context differs from the bacterial in many significant
29 ways. One key difference is that, in bacteria, regulatory DNA is found proximal to the gene, typically within
30 1kb upstream of the transcription start site (TSS), whereas eukaryotic regulatory sequences are found in
31 discrete pieces that may be proximal to, or distal from, the TSS. The eukaryotic regulatory elements known
32 as enhancers form a particularly important class. Enhancers were originally defined as DNA sequences which
33 could drive the expression of genes in a location and orientation independent manner [2, 3]. Since these initial
34 discoveries, many native enhancers have been identified which play critical roles in a variety of processes,
35 such as embryonic development [4], physiology [5] and evolution [6]. Genetic variation in enhancers has
36 also been shown to mediate risk for complex disease [7], Mendelian disease [8] and cancer [9]. Based on
37 these and many other studies, we know that enhancers can be located over 1Mb from a target gene TSS,
38 an individual enhancer may regulate multiple genes, some genes are regulated by multiple enhancers and
39 the set of enhancers actively regulating a given gene may depend on cellular context. These properties have
40 made it difficult to identify the rules governing enhancer-gene regulation.

41 Given their importance, much attention has been given to systematically identifying enhancer sequences
42 and the genes they regulate. An important breakthrough has been the development of high-throughput
43 CRISPR interference (CRISPRi) screens, which enable putative enhancer sequences to be perturbed in cell

44 culture and the resulting effect on expression of a target gene to be measured [10–13]. These screens typically
45 measure quantitative effects on gene expression as the proportional change in mean gene expression over a
46 cell population. We call this quantity the *fractional change* and, given its importance in this paper, define
47 it formally as follows: let $\psi(g)$ denote the wild-type mean expression level of a gene g , in whatever units
48 are used to measure it, and let $\psi(g, e_q)$ denote the mean expression level of g after an enhancer of g , e_q , has
49 been perturbed. The *fractional change* of e_q is then the non-dimensional quantity,

$$f(g, e_q) := \frac{\psi(g) - \psi(g, e_q)}{\psi(g)}. \quad (1)$$

50 The fractional change for thousands of putative enhancer-gene connections has been measured and compu-
51 tational methods have assessed whether the observed fractional change is statistically different from zero.
52 Current efforts are now focused on two main questions. First, what can we learn about enhancer biology
53 from these screens? Second, can the results from these screens be used to develop computational methods
54 which can predict which enhancers regulate which genes in arbitrary cellular contexts?

55 The Activity-by-Contact (ABC) model has been proposed as a way to make progress on both of these
56 questions [11]. The ABC model is based on the mechanistic notion that an enhancer’s effect on gene
57 expression depends on the intrinsic strength of the enhancer (activity) and the frequency with which it
58 comes into physical proximity to the gene promoter (contact). The ABC model gives rise to the ABC
59 Score, a quantitative formula which is intended to predict the fractional change observed in an enhancer
60 perturbation experiment. For a gene, g , with N putative enhancers, e_1, \dots, e_N , the ABC Score for a specific
61 enhancer e_q , is given by,

$$\text{ABC}(g, e_q) := \frac{\alpha_q \gamma_q}{\alpha_1 \gamma_1 + \dots + \alpha_N \gamma_N}, \quad (2)$$

62 where α_i represents the activity of e_i and γ_i represents the contact frequency between e_i and the promoter of
63 g . In [11] a putative enhancer was defined as a chromatin-accessible DNA element of approximately 500 base
64 pairs; α_i was assigned using measures of chromatin state of the enhancer, such as DNase-Seq and H3K27ac
65 ChIP-Seq; γ_i was assigned using the contact frequency between a putative enhancer and the gene promoter,
66 as measured by Hi-C; and the sum in the denominator of Eqn.2 was taken over all putative enhancers within
67 5Mb of g .

68 The ABC Score is reasonably effective at predicting the results of CRISPRi screens. When considered
69 as a binary classifier, the ABC Score has achieved a precision of 59% at 70% recall benchmarked against a
70 database of nearly 4,000 putative enhancer-gene connections in the K562 cell line [11]. Similar performance
71 has also been observed in other cell types [11, 14] and in subsequent benchmarking against other CRISPRi
72 screens in K562 cells [15, Fig.S8a]. We emphasize that the ABC Score is computed directly from genomic
73 data orthogonal to the CRISPRi experiment. As such, it has no free parameters and does not require
74 fitting or training. The classification ability of the ABC Score and its modest input data requirements have
75 resulted in its widespread use to interpret non-coding genetic variation [14, 16, 17], identify enhancers in
76 disease related contexts [18, 19] and investigate the dosage effect of transcription factor concentration on
77 gene expression [20].

78 Despite its practical utility as a binary classifier, the ability of the ABC Score to predict the fractional
79 change is fundamentally limited [11, Fig.3c]. From Eqn. 2, it is clear that the sum of the ABC Scores over
80 all putative enhancers of a given gene is equal to 1,

$$\sum_{i=1}^N \text{ABC}(g, e_i) = 1. \quad (3)$$

81 We define the total fractional change of a gene to be the sum of the fractional changes of all enhancers for
82 the gene, $f(g, e_1) + \dots + f(g, e_N)$. If the ABC model were perfectly reflecting the fractional change, so that
83 $\text{ABC}(g, e_i) = f(g, e_i)$, it would predict that the total fractional change for all genes is equal to 1. However,
84 experimentally, a range of total fractional changes has been observed from 0 to greater than 3 [10, 11, 14,
85 21–23]. This incompatibility is a consequence of the algebraic structure of the ABC Score formula and

86 cannot be resolved in a straightforward way. For example, it cannot be resolved by using different types of
87 epigenomic data to assign values to α_i or γ_i .

88 What, if anything, about enhancer biology can be concluded from the successes and limitations of the
89 ABC Score? We believe that considering this question requires a formal description of the ABC model. The
90 original description of the ABC model is *informal*, in the sense that the relationships between the mechanisms
91 of activity and contact and the ABC Score formula were not determined by formal mathematical arguments.
92 In consequence, the biological and biophysical assumptions that underlie formulas of this kind have not been
93 clarified.

94 In the present paper, we present a strategy for the formal mathematical modeling of enhancer-gene
95 regulation. We introduce the *default model*, a set of assumptions for how multiple enhancers independently
96 regulate a gene. We show that a formula with the same algebraic structure as the ABC Score formula in
97 Eqn.2 can be rigorously derived from a special case of the default model. This clarifies the assumptions that
98 underlie the ABC Score formula. However, these assumptions also imply that the total fractional change
99 of a gene is equal to one. We show how changing the assumptions of the default model can lead to total
100 fractional changes which are less than or greater than one. More generally, the framework introduced here
101 offers a rigorous foundation for future studies of enhancer-gene regulation.

102 Results

103 An activation-communication model of enhancer function

104 Our approach to modelling enhancer-gene regulation is based on the linear framework, a method of using
105 graphs to analyse biomolecular systems [24–26] that has been previously introduced to study gene regulation
106 [26]; see [27, 28] for up-to-date reviews. The graphs in question have vertices that are linked by labelled,
107 directed edges. The vertices represent molecular states of DNA, the edges represent transitions between these
108 molecular states and the labels represent the transition rates, which are positive numbers with dimensions
109 of (time)⁻¹.

110 An example linear framework graph is shown in Fig.1a. This graph, which we have called H , represents a
111 single enhancer which can be either activated (filled red circle) or not and in communication with its target
112 gene (curved arrow) or not. It thereby captures the two main notions in the original ABC model, although
113 we prefer to speak here of “communication” rather than of “contact” (see below). These two features of the
114 enhancer are treated in the graph as being independent of each other: the rates for becoming activated or
115 deactivated do not depend on the state of communication, and the rates for making or losing communication
116 do not depend on the state of activation. Independence will be one of the central features of our treatment
117 and will appear both in how an individual enhancer is treated, as in this example in Fig.1a, and in how a
118 gene is regulated by multiple enhancers, as we will explain below.

119 The graph H in Fig.1a represents a *coarse-graining* of the actual complexity of enhancer-gene regulation
120 (Fig.1b). Activation is intended to capture processes local to the enhancer sequence such as transcription
121 factor binding, chromatin reorganisation, nucleosome remodelling, recruitment of co-regulators or transcrip-
122 tion of the enhancer sequence itself to generate enhancer RNA. Communication refers to the processes by
123 which information is transferred from the enhancer to its target gene. Many communication mechanisms
124 have been proposed including physical contact through DNA looping [29], diffusion of regulatory molecules
125 [30] and phase separation [31]. We thus use the word ‘communication’ instead of ‘contact’ to reflect that
126 enhancer-gene regulation may not require physical contact. It is, of course, possible that the specific ac-
127 tivation or communication mechanisms may differ between enhancers. The value of this coarse-graining
128 lies in not making commitments about the underlying mechanism, at the price of ignoring the potential
129 consequences of how activation and communication are implemented in molecular terms. This particular
130 coarse-graining will facilitate our clarification of the ABC Score formula below.

131 Having provided an example of a linear framework graph and explained how it describes the biological
132 context that we will be studying, we now go into the details of the linear framework. We will make use of
133 the example in Fig.1a throughout this work.

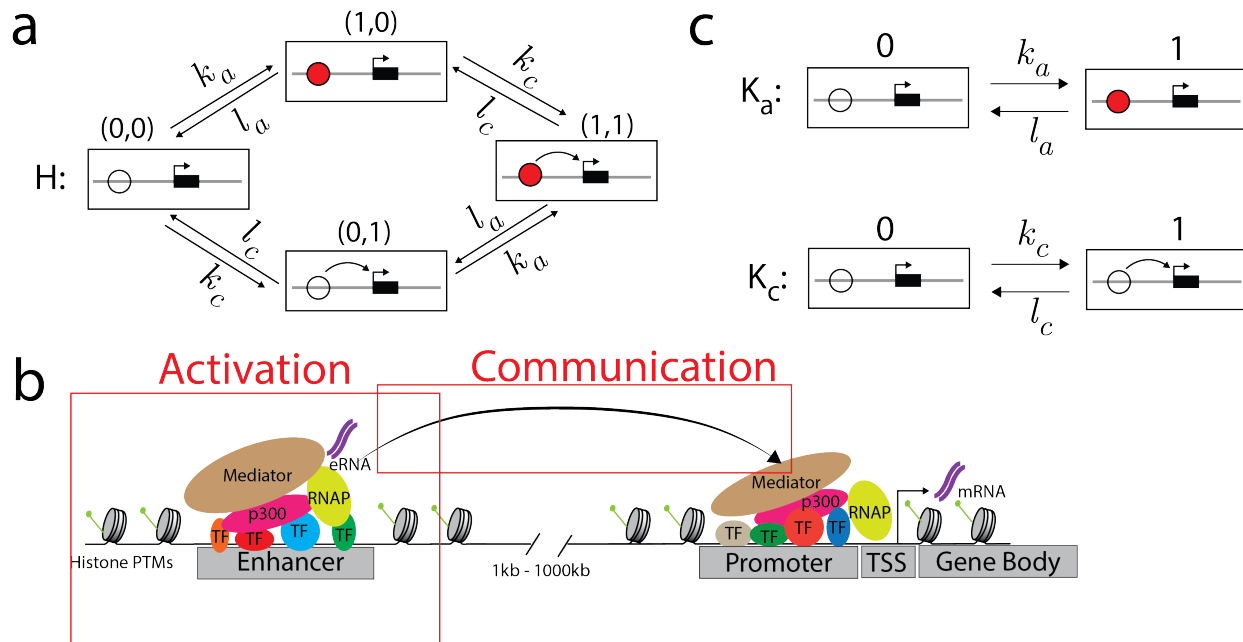


Figure 1: The activation-communication coarse graining. **a)** An example linear framework graph, H , representing a coarse-grained view of an enhancer. Each vertex contains a schematic of the enhancer (circle) and its target gene (black rectangle with the transcription start site marked with an arrow). The enhancer may be activated (filled red circle) or communicating (curved arrow to the target gene), encoded in the notation (i, j) used to denote vertices. The edge labels show that activation and communication take place independently of each other. **b)** A more detailed picture of the molecular complexity that may underlie the coarse-grained graph in panel **a**, as described further in the text. **c)** The example graph H in panel **a** is the graph product of two simpler 2-vertex graphs, K_a , which represents activation, and K_c , which represents communication. The product structure of H is equivalent to the independence of activation and communication.

134 Preliminaries on the linear framework

135 Notation and terminology

136 We will start by introducing some basic ideas about linear framework graphs. We will use a letter like G
 137 or H to refer to a graph. Vertices will generally be denoted i, j , etc. We will use the notation $i \in G$ to
 138 mean the state i from the graph G . Edges will be denoted $i \rightarrow j$ and edge labels will be denoted $\ell(i \rightarrow j)$.
 139 So, using the notation for the example graph H in Fig.1a, $\ell((0, 0) \rightarrow (0, 1)) = \ell((1, 0) \rightarrow (1, 1)) = k_c$ (the
 140 notation for the vertices in this graph arises from its product structure and will be explained later). If some
 141 feature X is being discussed for different graphs, we will sometimes use brackets, as in $X(G)$, or a subscript,
 142 as in X_G , to specify the graph in question. We will use the word *structure* to refer to just the vertices and
 143 edges of a graph, ignoring the edge labels; when we say “graph”, we will always be including the labels, even
 144 when they are not mentioned explicitly.

145 The Markov process

146 A graph G is equivalent to a finite-state, continuous-time, time-homogeneous Markov process [25, 28, 32].
 147 This stochastic behaviour can be understood as follows. If the system is in state i , then for each edge $i \rightarrow j$
 148 which leaves i , a “firing” time is randomly chosen from the exponential probability distribution, $\lambda \exp(-\lambda t)$,
 149 where λ is the transition rate of that edge, $\lambda = \ell(i \rightarrow j)$, and the edge with the lowest firing time is taken,

150 at that time. This generates a stochastic trajectory of states and transitions. If we follow a trajectory up
 151 to time T and measure the proportion of time spent in state i , then that ratio stabilises with increasing T
 152 to become the *steady-state probability* of state i [32], which we will denote by $u_i^*(G)$. Provided G is *strongly*
 153 *connected*, this quantity does not depend on the state in which the trajectory starts and the steady-state
 154 probability is a property of the graph [25]. A strongly connected graph is one in which any two distinct
 155 vertices, i and $j \neq i$, are connected by a directed path, $i = i_1 \rightarrow i_2 \rightarrow \dots \rightarrow i_k = j$. The example graph H
 156 in Fig.1a is strongly connected but ceases to be if the edges $(1, 1) \rightarrow (1, 0)$ and $(1, 1) \rightarrow (0, 1)$ are removed.
 157 We will assume from now on that all our graphs are strongly connected.

158 Thermodynamic equilibrium and steady-state probabilities

159 One of the advantages of the linear framework is that, provided the graph is finite, its steady-state probabilit-
 160 ies can be calculated algebraically in terms of the edge labels. (We will encounter an infinite graph below
 161 but, as we will see, we do not have to deal with them directly and can work only with finite graphs.) If the
 162 graph can reach *thermodynamic equilibrium* the algebra can be done quite easily but, importantly, it can
 163 also be done when the graph is away from thermodynamic equilibrium, although the formulas become more
 164 complicated. A graph can reach thermodynamic equilibrium if, and only if, it satisfies two conditions. First,
 165 it must be *reversible*, so that if there is an edge $i \rightarrow j$, then there is also an edge $j \rightarrow i$, which represents the
 166 reverse of the process that corresponds to $i \rightarrow j$. Second, it must satisfy the *cycle condition*: the product of
 167 the label ratios around any cycle of reversible edges must be 1. The graph in Fig.1a is evidently reversible
 168 and has only one cycle of reversible edges, $(0, 0) \rightleftharpoons (1, 0) \rightleftharpoons (1, 1) \rightleftharpoons (0, 1) \rightleftharpoons (0, 0)$, for which the product
 169 of label ratios is

$$\left(\frac{k_a}{l_a}\right) \left(\frac{k_c}{l_c}\right) \left(\frac{l_a}{k_a}\right) \left(\frac{l_c}{k_c}\right) = 1. \quad (4)$$

170 For this graph, the independence of activation and communication ensures that the graph can reach ther-
 171 modynamic equilibrium.

172 When a graph can reach thermodynamic equilibrium, its steady-state probabilities can be calculated as
 173 follows. First, choose any vertex as a reference; let us call it 1. Second, choose any path of reversible edges
 174 from 1 to the state in question, say i : $1 \rightleftharpoons i_1 \rightleftharpoons \dots \rightleftharpoons i_k = i$. The steady-state probability of i is then
 175 proportional to the product of the label ratios along this path,

$$u_i^*(G) \propto \left(\frac{\ell(i_1 \rightarrow i_2)}{\ell(i_2 \rightarrow i_1)}\right) \times \dots \times \left(\frac{\ell(i_{k-1} \rightarrow i_k)}{\ell(i_k \rightarrow i_{k-1})}\right). \quad (5)$$

176 It is a simple consequence of the cycle condition that the quantity on the right-hand side of Eqn.5 does
 177 not depend on the choice of path from 1 to i . The proportionality constant in Eqn.5 is readily obtained
 178 by exploiting the fact that the sum of all the probabilities must be 1, so that, if the vertices are denoted
 179 $1, \dots, N$, then $u_1^*(G) + \dots + u_N^*(G) = 1$. If we follow this prescription for the graph in Fig.1a, we find that,
 180 for example,

$$u_{(1,1)}^*(G) = \frac{(k_a/l_a)(k_c/l_c)}{1 + (k_a/l_a) + (k_c/l_c) + (k_a/l_a)(k_c/l_c)}. \quad (6)$$

$$= \frac{k_a}{k_a + l_a} \cdot \frac{k_c}{k_c + l_c} \quad (7)$$

181 (We will sometimes use a "·" to denote multiplication to make formulas like this look clearer.) The reor-
 182 ganisation of Eqn.6 into Eqn.7 reveals a product structure in the algebra whose significance will emerge
 183 below. The formula in Eqn.6 is the same as would arise from equilibrium statistical mechanics. It is one
 184 of the features of the linear framework that it reduces to equilibrium statistical mechanics for systems that
 185 are at thermodynamic equilibrium but also yields algebraic formulas for systems away from thermodynamic
 186 equilibrium.

187 Product graphs as models of independence

188 In studying gene regulation, a very helpful construction is that of a *product graph*, because it captures the
 189 default situation in which two or more genetic systems operate independently of each other. The example
 190 graph in Fig.1a is a case in point. This graph H is the product of the graphs K_a and K_c in Fig.1c. Here,
 191 K_a is a two-vertex graph that represents just the activation of the enhancer and K_c is a two-vertex graph
 192 that represents just the communication.

193 We will use K_a and K_c to describe the product graph construction. We will do this in two steps. We
 194 will first specify the vertices and edges by building the *product structure*, denoted $K_a \times K_c$, and then we will
 195 specify the labels to get the *product graph*, denoted $K_a \otimes K_c$. As we will see below, product structures underlie
 196 other constructions in which the independence of the product graph is broken, which is why it is helpful to
 197 distinguish structures and graphs. The vertices in $K_a \times K_c$ are ordered pairs, (i, j) , of vertices $i \in K_a$ and
 198 $j \in K_c$. The edges in $K_a \times K_c$ arise from the edges in either component K_a or K_c , taken independently of
 199 the state of the other component. In other words, if $i_1 \rightarrow i_2$ is any edge in K_a , then $(i_1, j) \rightarrow (i_2, j)$ is an
 200 edge in $K_a \times K_c$, for all $j \in K_c$; similarly, if $j_1 \rightarrow j_2$ is any edge in K_c , then $(i, j_1) \rightarrow (i, j_2)$ is an edge in
 201 $K_a \times K_c$, for all $i \in K_a$; these are the only edges in $K_a \times K_c$. This prescription yields the structure of the
 202 graph H in Fig.1a.

203 The labels of the product graph, $K_a \otimes K_c$, are also inherited from those in K_a or K_c , independently of
 204 the state of the other component,

$$\ell_{K_a \otimes K_c}((i_1, j) \rightarrow (i_2, j)) = \ell_{K_a}(i_1 \rightarrow i_2) \quad \text{and} \quad \ell_{K_a \otimes K_c}((i, j_1) \rightarrow (i, j_2)) = \ell_{K_c}(j_1 \rightarrow j_2).$$

205 We see that $K_a \otimes K_c$ corresponds exactly to the graph H in Fig.1a. The graph product precisely captures
 206 the sense in which the components of the product, here K_a and K_c , operate independently of each other:
 207 the transitions in either component are unaffected, as to their occurrence and their rates, by the state of the
 208 other component.

209 In the more general case of a product of m graphs, the vertices are naturally indexed as ordered tuples,
 210 (i_1, \dots, i_m) .

211 One of the consequences of the product graph construction is that its steady-state probabilities are easily
 212 calculated. If K_1, \dots, K_N are any set of N strongly connected graphs, then the steady-state probabilities
 213 in the product graph $K_1 \otimes \dots \otimes K_N$ can be computed by multiplying the steady-state probabilities in the
 214 individual graphs,

$$u_{(i_1, \dots, i_N)}^*(K_1 \otimes \dots \otimes K_N) = u_{i_1}^*(K_1) \cdots u_{i_N}^*(K_N). \quad (8)$$

215 Eqn.8, which is proved in [26], again captures the sense in which the components K_1, \dots, K_N are indepen-
 216 dent of each other. We note that Eqn.8 holds even for graphs which are unable to reach thermodynamic
 217 equilibrium.

218 We can see Eqn.8 at work for the graph in Fig.1a, which is the product of the graphs K_a and K_c in
 219 Fig.1c. If we follow the prescription in Eqn.5, we see that

$$u_1^*(K_a) = \frac{k_a/l_a}{1 + k_a/l_a} \quad \text{and} \quad u_1^*(K_c) = \frac{k_c/l_c}{1 + k_c/l_c}. \quad (9)$$

220 If we apply Eqn.8 to the formulas above, we see that,

$$\begin{aligned} u_{(1,1)}^*(K_a \otimes K_c) &= \left(\frac{k_a/l_a}{1 + k_a/l_a} \right) \left(\frac{k_c/l_c}{1 + k_c/l_c} \right) \\ &= \frac{k_a}{k_a + l_a} \cdot \frac{k_c}{k_c + l_c}, \end{aligned} \quad (10)$$

221 which recovers the expression in Eqn.7, whose algebraic product structure is now seen to reflect the underlying
 222 product graph.

223 Eqn.8 for individual vertices has a straightforward extension to subsets of vertices. To explain this, let
 224 K be any graph and let $S \subseteq K$ be any subset of vertices in K . The steady-state probability of being in any

225 vertex of S , denoted $u_S^*(K)$, is given by $u_S^*(K) = \sum_{i \in S} u_i^*(K)$. Now suppose, as above, that K_1, \dots, K_N
 226 are any strongly connected graphs. Let $S_i \subseteq K_i$ be any subset of vertices of K_i and let $S_1 \times \dots \times S_N$ be
 227 the corresponding *set product* in K . This set product, for which we use, for convenience, the same notation
 228 as for the product structure, has the obvious definition that it consists of all those tuples (i_1, \dots, i_N) where
 229 $i_k \in S_k$. It is then a simple consequence of Eqn.8 that,

$$u_{S_1 \times \dots \times S_N}^*(K_1 \otimes \dots \otimes K_N) = u_{S_1}^*(K_1) \dots u_{S_N}^*(K_N). \quad (11)$$

230 One of the implications of Eqn.11 is that if we take i_j to be a coordinate that runs over the vertices of K_j ,
 231 then the probability that i_j has a particular value, say $i_j = b$, remains the same irrespective of the other
 232 factors in the graph product,

$$u_{\{i_j=b\}}^*(K_1 \otimes \dots \otimes K_N) = u_{\{i_j=b\}}^*(K_j). \quad (12)$$

233 Eqn.12 follows from Eqn.11 because the subset $\{i_j = b\}$ in $K_1 \otimes \dots \otimes K_N$ is the product subset,

$$K_1 \times \dots \times K_{j-1} \times \{i_j = b\} \times K_{j+1} \times \dots \times K_N,$$

234 and $u_{K_i}^*(K_i) = 1$. We can see an example of Eqn.12 at work in Figure 1. Let $\{a=1\} = \{(1,0), (1,1)\}$ be the
 235 subset of vertices of H in which the enhancer is activated. Then Eqn.12 shows that $u_{\{a=1\}}^*(H) = u_1^*(K_a)$.
 236 We will make further use of Eqn.12 in what follows.

237 The gene expression response

238 The graphs we have considered up to now are models of the regulatory state of the gene. We now discuss
 239 how to incorporate the production and degradation of mRNA. The standard approach in the literature is
 240 known as *kinetic modeling* and uses a Markovian framework based on the chemical master equation [33]. We
 241 follow this same approach within the graph-theoretic setting introduced here.

242 At any given time, the state of gene expression is specified by a certain number of molecules of the
 243 corresponding mRNA. This number increases by 1 each time RNA polymerase transcribes the gene and
 244 decreases by 1 each time an mRNA molecule is degraded or lost through transport out of the nucleus. We
 245 can represent such an expression system by the (semi)-infinite *pipeline* structure, P , in which the state p
 246 represents the number p of mRNA molecules, from $p = 0$ onwards, and the edges correspond to mRNA
 247 production, $p \rightarrow p + 1$, and degradation or loss, $p \rightarrow p - 1$ (Fig.2a).

248 Given a gene-regulatory graph, G , we represent the overall system of regulation and expression by a
 249 *copy-number graph*, $G \times P$, that will be derived from the product structure, $G \times P$ (Fig.2b). The states of
 250 $G \times P$ are identical to those of $G \times P$ but $G \times P$ may not have all the edges that are present in $G \times P$.
 251 Each state in $G \times P$ keeps track of the regulatory state of the gene and the number of mRNA molecules
 252 that are present. We now discuss how to assign labels to this graph (Fig.2b). We assume that each state,
 253 $i \in G$, has a corresponding non-negative rate of mRNA production, $r_i(G) \geq 0$. If $r_k(G) = 0$, so that mRNA
 254 production is not possible in state k of G , then the edges $(k, p) \rightarrow (k, p + 1)$ are removed from $G \times P$ for all
 255 $p \in P$. (Note that edge labels must always be positive.) This is the only way in which the structure of $G \times P$
 256 differs from that of $G \times P$. If $r_k(G) > 0$, we will assume that the rate of mRNA production does not change
 257 with the number of mRNA molecules that have been expressed, so that $\ell((k, p) \rightarrow (k, p + 1)) = r_k(G)$ for
 258 any $p \in P$. As for mRNA degradation or loss, this takes place independently of the regulatory system, so
 259 the most parsimonious assumption is that its rate is proportional to the number of mRNAs that are present
 260 and is independent of the regulatory state. Accordingly, we may write $\ell((k, p) \rightarrow (k, p - 1)) = \delta(G) \cdot p$ for
 261 any $k \in G$ and any positive $p \in P$, where $\delta(G)$ is the degradation rate constant. Finally, we assume that
 262 regulatory transitions do not depend on gene expression, so that $\ell((i, p) \rightarrow (j, p)) = \ell_G(i \rightarrow j)$ for all $i, j \in G$
 263 and for all $p \in P$. A compact way to visually represent a copy-number graph is shown in Fig.2c.

264 At steady state, $G \times P$ gives rise to a probability distribution over the mRNA copy number. We will
 265 define the response of the gene, which we will denote by $R(G)$, to be given by the average of this number
 266 distribution,

$$R(G) := \sum_{(i,p)} p \cdot u_{(i,p)}^*(G \times P). \quad (13)$$

267 We note that $R(G) \geq 0$.

268 Because $G \times P$ is not a finite graph, the prescription given in Eqn.5 for calculating steady-state proba-
269 bilities no longer works. ($G \times P$ is also not at thermodynamic equilibrium, unless every regulatory state has
270 the same rate of mRNA production, as can be checked by following the cycle condition formula in Eqn.4.)
271 However, we can appeal to a very useful theorem, due to Sanchez and Kondev, which tells us that we do
272 not have to operate on $G \times P$ in order to calculate $R(G)$ [34]. Translating their work into the graph-theory
273 language used here, we find that the response of the gene can be calculated in terms of the average of $r_i(G)$
274 over only the the steady-state probabilities of G , normalized by $\delta(G)$,

$$R(G) = \frac{1}{\delta(G)} \sum_{i \in G} r_i(G) \cdot u_i^*(G). \quad (14)$$

275 It follows that, although infinite graphs arise to represent the mRNA expression system, we do not need to
276 work with them to calculate the mean steady-state expression $R(G)$, under the assumptions made above.
277 Sanchez and Kondev did not use graph theory in their work, so we provide an independent graph-based
278 proof of Eqn.14 in the Methods. In subsequent work, we will show how the copy-number graphs introduced
279 here lead to generalizations of the results of [34] but we do not need that for the present paper.

280 This result provides some justification for reducing the notational clutter from multiple instances of P .
281 We will refer to the *regulatory graph* G when we mean G on its own, and to the *copy-number graph* G when
282 we mean $G \times P$, defined for some specified choice of production rates $r_i(G)$ and degradation rate $\delta(G)$. These
283 parameters may not be explicitly mentioned when speaking of a copy-number graph but they should be kept
284 in mind.

285 As an illustration of Eqn.14, we will assign production rates to the graph in Fig.1a and compute its
286 response. We will make the assumption that mRNA is only produced when the enhancer is both activated
287 and communicating (Fig. 2d). The mRNA production rates of H are therefore given by

$$r_{(0,0)}(H) = r_{(1,0)}(H) = r_{(0,1)}(H) = 0 \quad \text{and} \quad r_{(1,1)}(H) = r. \quad (15)$$

288 We can now use Eqn.14 to calculate the response, $R(H)$, taking advantage of Eqn.10, in which we exploited
289 the product graph decomposition $H = K_a \otimes K_c$. We see that,

$$R(H) = \frac{r}{\delta} \cdot \left(\frac{k_a}{k_a + l_a} \right) \left(\frac{k_c}{k_c + l_c} \right). \quad (16)$$

290 It follows from Eqns.9 and 12 that we can interpret $k_a/(k_a + l_a)$ as the probability that the enhancer is
291 activated and, similarly, $k_c/(k_c + l_c)$ as the probability that the enhancer is communicating. Eqn.16 tells us
292 that the response of H is the product of the ratio of production to degradation, the probability of activation
293 and the probability of communication.

294 This concludes our analysis of a gene regulated by a single enhancer using the activation-communication
295 coarse graining. We now turn to considering how multiple enhancers work together to regulate gene expres-
296 sion.

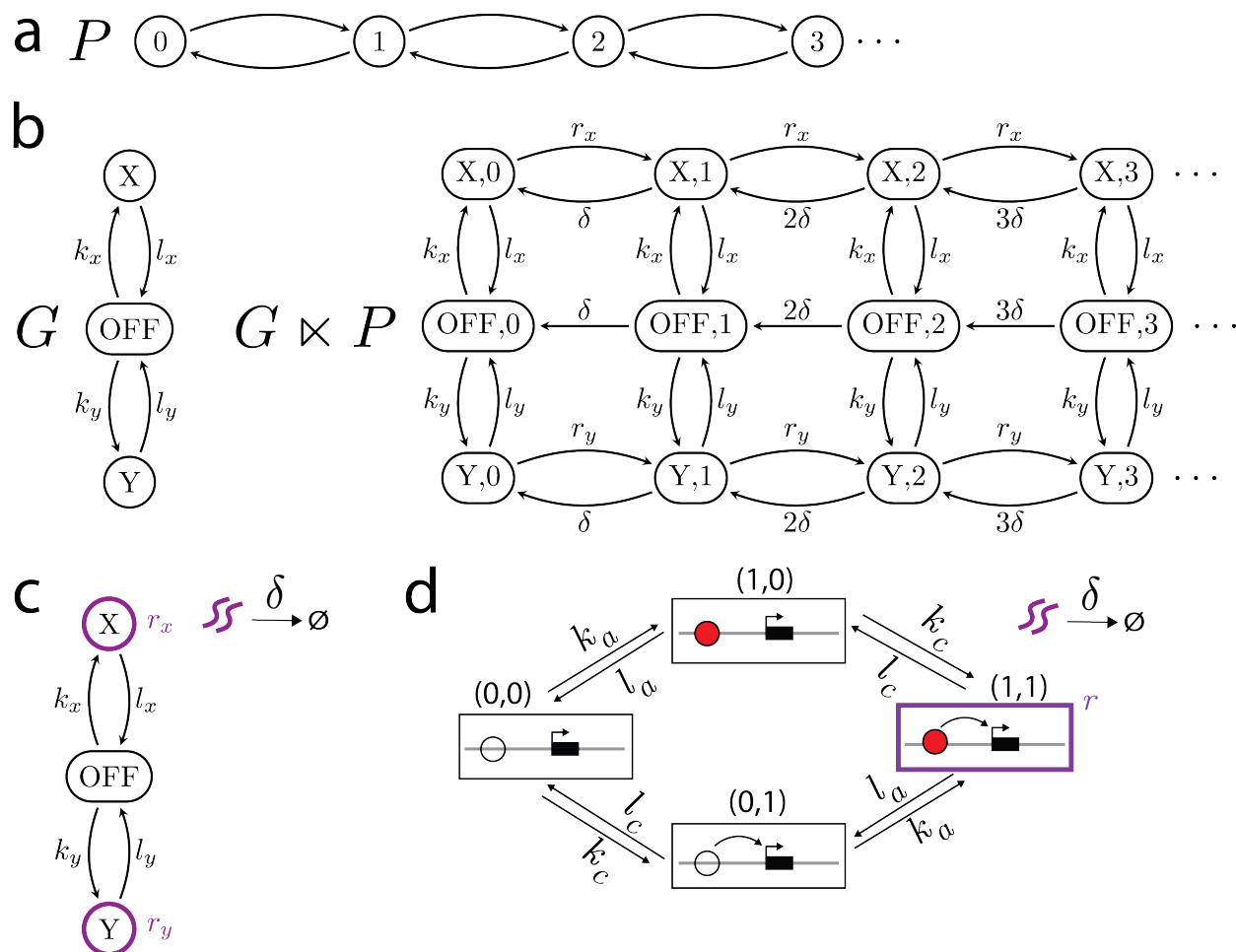


Figure 2: Modeling mRNA production and degradation through copy-number graphs. **a)** The pipeline structure P represents the number of mRNA molecules and their production and loss. **b)** An example regulatory graph, G , and the resulting copy-number graph $G \times P$. In this example G has two production states, X and Y , with corresponding mRNA production rates r_x and r_y respectively. We note that $G \times P$ is a sub-structure of $G \times P$; it has the same vertices but lacks the edges corresponding to a production rate of zero. G and P also do not operate independently in $G \times P$ because the mRNA production rates depend on the regulatory state. **c)** A compact way to represent $G \times P$. The production states are outlined in purple with corresponding mRNA production rates. The degradation rate, δ , is shown above the arrow from purple squiggles (mRNA) to the empty set \emptyset . **d)** The compact representation of the graph $H \times P$, where H is given in Fig.1a. The only production state is the state $(1, 1)$, in which the enhancer is both activated and communicating, which has production rate r .

297 A default model of how multiple enhancers independently regulate a gene

298 We now introduce the *default model* of enhancer-gene regulation. This is a set of assumptions for how
 299 multiple enhancers collectively regulate a gene in an independent manner. We previously introduced the
 300 product graph construction which represents independence between regulatory graphs. We now broaden
 301 those assumptions to also allow for mRNA production. We expect this default model construction to be of
 302 general interest. In the next section we will show how a special case of the default model clarifies the ABC
 303 Score formula.

304 Consider a gene, g , that is regulated by N enhancers, e_1, \dots, e_N . We will assume that enhancer e_l is
305 modelled by the graph G_l . We make no assumptions about G_l other than the prevailing assumption that
306 all our graphs are strongly connected. G_l could be substantially more complicated than the graph in Fig.1a
307 and could incorporate, for example, chromatin organisation, nucleosomes, co-regulators, post-translational
308 modifications, chromosome conformation, etc [26]. In particular, there is no requirement that G_l should be
309 able to reach thermodynamic equilibrium. At this point our assumptions are very general and could apply
310 to essentially any enhancer, when considered from a Markovian perspective.

311 We denote the graph that models the collective regulation of the enhancers by G and describe how G is
312 defined in terms of the G_l .

313 The first assumption says that each enhancer has its own individual effect.

314 1. Individuality. Each enhancer e_l , when acting in the absence of any of the other enhancers, drives gene
315 expression at the rate $r_i(G_l) \geq 0$ for each state $i \in G_l$, and gives rise to the response $R(G_l)$, as defined
316 by Eqn.13. If the enhancer is unable to drive expression on its own, then $r_i(G_l) = 0$ for every state
317 $i \in G_l$.

318 The next two assumptions specify how the enhancers work together.

319 2. Regulatory independence. Each enhancer acts independently of all the others, so that the regulatory
320 graph of G is given by the product graph $G_1 \otimes \dots \otimes G_N$.

321 3. Production-rate summation. Each enhancer independently influences mRNA production. Accordingly,
322 if (i_1, \dots, i_N) is a state in G , then its mRNA production rate is a sum of the corresponding production
323 rates in each enhancer graph:

$$r_{(i_1, \dots, i_N)}(G) = r_{i_1}(G_1) + \dots + r_{i_N}(G_N). \quad (17)$$

324 The summation of rates in Assumption 3 arises for the following reason. If each enhancer influences mRNA
325 production independently, then the time at which an mRNA is produced will be the minimum of the times
326 at which each individual enhancer has its effect on production. These individual times are exponentially
327 distributed with rates $r_{i_j}(G_j)$ for state i_j in G_j . The minimum of several exponentially distributed random
328 variables is a random variable that is also exponentially distributed, with rate given by the sum of the
329 individual rates. This leads to Eqn.17. The final assumption specifies the degradation rates.

330 4. Uniform degradation. Since mRNA degradation is a separate process to gene regulation and gene
331 expression, we consider the characteristic degradation rate to be a property of the gene, not the
332 enhancer. As such, each graph G_l is assumed to have the same degradation rate, $\delta(G_l) = \delta$ for all l ,
333 and the mRNA degradation rate of G is also δ : $\delta(G) = \delta$.

334 For any set of copy-number graphs G_1, \dots, G_N , we denote the copy-number graph which models their
335 collective effect on transcription according to Assumptions 1 to 4 by the graph

$$G_1 \otimes \dots \otimes G_N. \quad (18)$$

336 Example constructions using the default model are given in Fig.3.

337 Assumptions 1 to 4 specify our default model of how enhancers collectively regulate a gene. Whether
338 any of the default model assumptions hold for an individual gene is a question that has to be addressed
339 experimentally. In particular, we would expect that Assumption 3 would eventually break down as more
340 enhancers are added to a gene since production rates will be limited by the physical processes involved in
341 transcription. Our goal here is to rigorously work out the consequences of these assumptions, so that we
342 know what to expect when the assumptions do hold and can compare these predictions to what is found
343 experimentally. Of particular significance is that the assumptions above imply that the collective response
344 of the enhancers is always the sum of their individual responses,

$$R(G_1 \otimes \dots \otimes G_N) = R(G_1) + \dots + R(G_N). \quad (19)$$

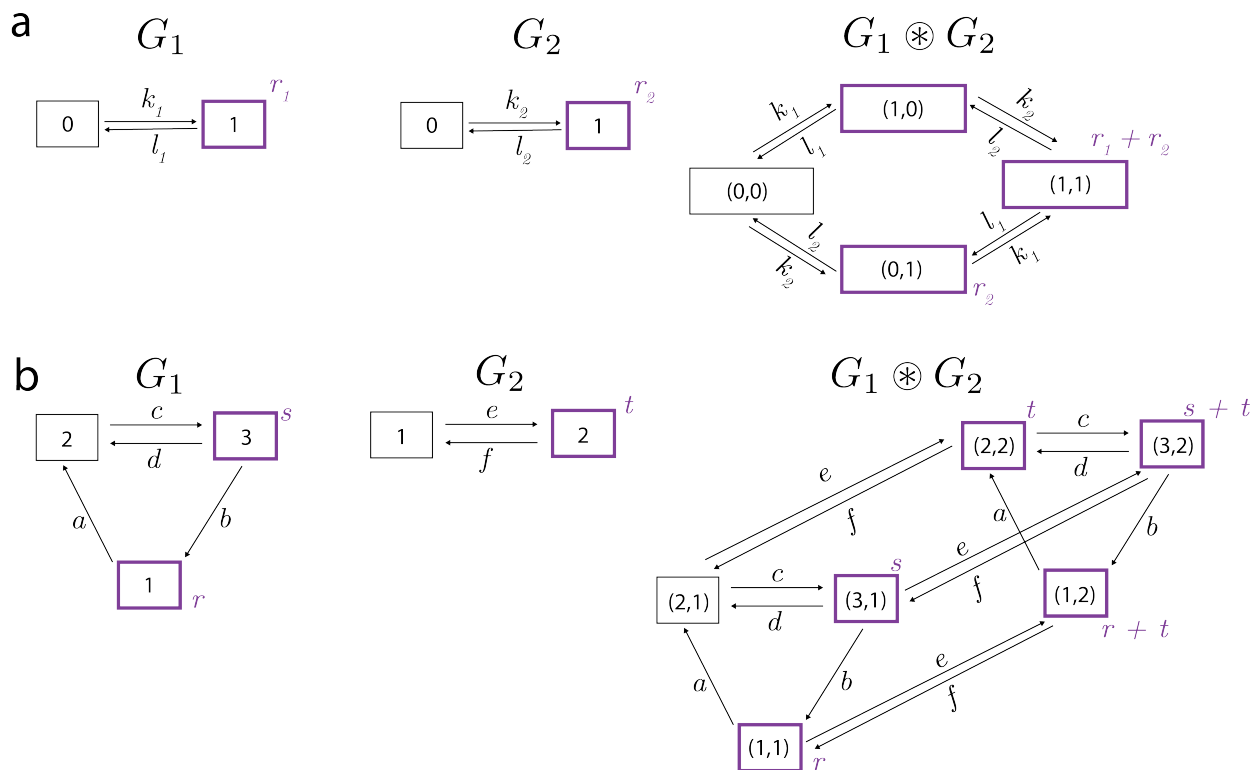


Figure 3: Two examples **a** and **b** of the default model construction. The copy-number graphs are depicted in compact format, as shown in Fig.2c but omitting the degradation symbols for clarity. Production states are outlined in bold purple with corresponding production rates in purple text. The model in **b** is adapted from Figure 9 of [26].

345 A proof of this fundamental property of the default model is given in the Methods.

346 A recent commentary has argued that formal definitions and rigorous modeling are necessary to investigate whether a set of enhancers is “greater than the sum of its parts” [35]. We fully agree and suggest
 347 that the notion of independence encoded by the default model, which gives rise to Eqn.19, could serve as a
 348 definition of what it means for a gene to be the *sum of its parts*.
 349

350 Transcription in the default model relies on the presence of enhancers. It is well known that the promoter
 351 sequences at some eukaryotic genes are sufficient to drive transcription even in the absence of distal enhancers
 352 [36, 37]. It is a future area of research to incorporate the role of core promoter elements and promoter proximal
 353 regulatory sequences along with their interactions with distal enhancers.

354 A clarification of the ABC Score formula

355 Enhancer perturbation and deletion fidelity

356 To see how formulas similar to the ABC Score can be derived from the default model, we need to consider
 357 how to formally model perturbations to enhancers such as genetic deletions or CRISPRi. As previously,
 358 we will assume that the target gene g is collectively regulated by enhancers e_1, \dots, e_N . We assume that
 359 enhancer e_i is modeled by the graph G_i and that g is modeled by $G_1 \otimes \dots \otimes G_N$. Let us consider what
 360 happens when enhancers e_{l_1}, \dots, e_{l_k} are perturbed in such a way that they are considered to no longer be
 361 working to regulate g . We will use a similar notation to that for the fractional change in the Introduction and
 362 denote the graph that arises from this perturbation as $G|_{e_{l_1}, \dots, e_{l_k}}$. In analogy to the fractional change,

363 we can define the *deletion effect* of the perturbation, $\Delta(G; e_{l_1}, \dots, e_{l_k})$, to be the proportional change in
 364 response of g ,

$$\Delta(G; e_{l_1}, \dots, e_{l_k}) := \frac{R(G) - R(G|e_{l_1}, \dots, e_{l_k})}{R(G)}. \quad (20)$$

365 It is important to keep in mind that the deletion effect is defined in terms of a model of gene regulation,
 366 whereas the fractional change is defined in terms of experimental data. The definition in Eqn.20 implicitly
 367 assumes that the system has returned to steady state following the perturbation. Furthermore, Eqn.20 says
 368 nothing about how the enhancer is perturbed or whether a CRISPRi perturbation has the same effect as a
 369 genetic deletion.

370 To calculate the deletion effect using Eqn.20, we need to know the perturbed graph, $G|e_{l_1}, \dots, e_{l_k}$. We
 371 assume that the graph shows *deletion fidelity*, which implies that the perturbation completely abrogates the
 372 function of the targeted enhancers and does not influence other enhancers. Let m_1, \dots, m_p be the remaining
 373 indices in $1, \dots, N$ after l_1, \dots, l_k have been removed.

374 5. Deletion fidelity. The regulatory graph of $G|e_{l_1}, \dots, e_{l_k}$ is the graph product $G_{m_1} \otimes \dots \otimes G_{m_p}$, and the
 375 production rates in $G|e_{l_1}, \dots, e_{l_k}$ are directly inherited from G ,

$$r_{(i_{m_1}, \dots, i_{m_p})}(G|e_{l_1}, \dots, e_{l_k}) = r_{i_{m_1}}(G_{m_1}) + \dots + r_{i_{m_p}}(G_{m_p}). \quad (21)$$

376 Deletion fidelity ensures that if G obeys Assumptions 1-4, then $G|e_{l_1}, \dots, e_{l_k}$ also obeys Assumptions 1-4
 377 for the remaining enhancers e_{m_1}, \dots, e_{m_p} and that, for the copy-number graphs,

$$G|e_{l_1}, \dots, e_{l_k} = G_{m_1} \circledast \dots \circledast G_{m_p}. \quad (22)$$

378 Using Eqn.22, it follows from Eqn.19 that,

$$R(G|e_{l_1}, \dots, e_{l_k}) = R(G_{m_1}) + \dots + R(G_{m_p}), \quad (23)$$

379 and so the formula for the deletion effect in Eqn.20 tells us that,

$$\Delta(G; e_{l_1}, \dots, e_{l_k}) = \frac{R(G_{l_1}) + \dots + R(G_{l_k})}{R(G_1) + \dots + R(G_N)}. \quad (24)$$

380 Eqns.23 and 24 are general properties that hold for the default model whenever Assumption 5 of deletion
 381 fidelity also holds. They allow us to formalise the notion of enhancer *additivity*, which we will discuss below,
 382 but, first, let us turn to the ABC Score formula.

383 The Independent-Activation-Communication (IAC) model

384 In the default model, the graph representing each individual enhancer can be arbitrarily complicated. To
 385 show how the ABC Score formula can arise from the default model, we need to impose the further assumption
 386 that each enhancer is modeled by the activation-communication coarse graining shown in Figs.1a and 1b.

387 6. The activation-communication coarse-graining. Enhancer e_i is described by the graph H_i , where H_i is
 388 the same graph as H in Fig.2d. Specifically, H_i is the graph product of an activation graph, $K_{a,i}$, with
 389 labels $k_{a,i}, l_{a,i}$, and a communication graph, $K_{c,i}$, with labels $k_{c,i}, l_{c,i}$ (Fig.1c), and $H_i = K_{a,i} \otimes K_{c,i}$.
 390 The only non-zero production rate of H_i occurs in the state in which the enhancer is both active and
 391 communicating, where the rate is r_i .

392 The overall regulatory system is then described by $G = H_1 \circledast \dots \circledast H_N$ (Fig. 4, Fig.S1). We call the model
 393 obeying Assumptions 1-6 the Independent-Activation-Communication (IAC) model. It follows from Eqn.16
 394 that the response of enhancer i in the IAC model is given by

$$R(H_i) = \frac{r_i}{\delta} \left(\frac{k_{a,i}}{k_{a,i} + l_{a,i}} \right) \left(\frac{k_{c,i}}{k_{c,i} + l_{c,i}} \right). \quad (25)$$

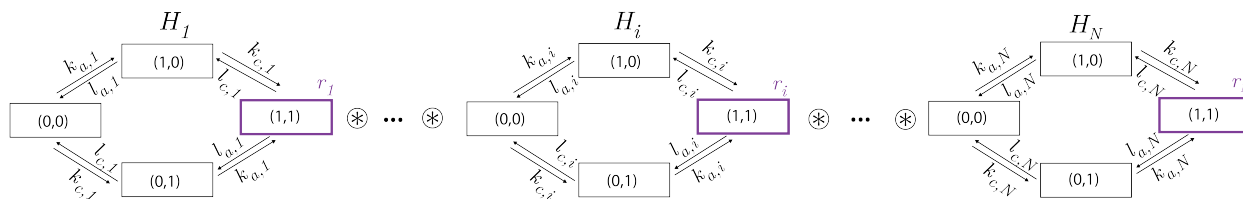


Figure 4: The Independent-Activation-Communication (IAC) model. A gene described by the IAC model follows Assumptions 1-6 for the component graphs H_1, \dots, H_N . The ordered pair of binary digits for the vertices in each H_i represent the activation and communication status, respectively, of each enhancer. Each H_i has the same structure as the graph in Fig.2d, but different labels, and represents the independence of activation and communication within each enhancer. Each H_i is assumed to have the same degradation rate which is omitted for clarity. See also Fig.S1.

395 Let us define $\tilde{\alpha}_i := k_{a,i}/(k_{a,i} + l_{a,i})$ and $\tilde{\gamma}_i := k_{c,i}/(k_{c,i} + l_{c,i})$ and recall from Eqn.16 that these quantities
 396 are the probability of activation and the probability of communication, respectively, of enhancer e_i . According
 397 to the fundamental property of the default model in Eqn.19, the response, $R(G)$, of the overall graph,
 398 $G = H_1 \otimes \dots \otimes H_N$, is given by,

$$R(G) = \frac{1}{\delta} \sum_{i=1}^N r_i \tilde{\alpha}_i \tilde{\gamma}_i. \quad (26)$$

399 Furthermore, as a consequence of deletion fidelity (Assumption 5), it follows from Eqn.24 that the deletion
 400 effect for enhancer e_q is given by,

$$\Delta(G; e_q) = \frac{r_q \tilde{\alpha}_q \tilde{\gamma}_q}{r_1 \tilde{\alpha}_1 \tilde{\gamma}_1 + \dots + r_N \tilde{\alpha}_N \tilde{\gamma}_N}. \quad (27)$$

401 Eqn.27 shows a striking algebraic similarity to the ABC Score formula in Eqn.2. The quantity $\tilde{\gamma}_i$, which
 402 is the probability of communication, is analogous to the ‘frequency of contact’, γ_i , that was envisaged
 403 for the ABC model [11] and appears in Eqn.2. There are different possible interpretations for the other
 404 terms. One potential interpretation for the term $r_i \tilde{\alpha}_i$ in Eqn.27, which is the production rate multiplied
 405 by the probability of activation, is that it is analogous to the ‘strength of the enhancer’, α_i , appearing in
 406 Eqn.2. Another potential interpretation is that $\tilde{\alpha}_i$ corresponds to α_i and that the production rates r_i are
 407 not represented in the ABC model. If the production rates are assumed to be equal, they would cancel
 408 out in Eqn.2, which would be consistent with a correspondence between $\tilde{\alpha}_i$ and α_i . Such interpretational
 409 ambiguities are to be expected because the ABC model is informal, while the IAC model presented here is
 410 formal. Moreover, our formal model separately specifies the regulatory state of the enhancer and its effect
 411 on transcription, whereas the ABC model does not make this distinction. An interesting question arises as
 412 to how numerical values can be assigned to the terms in Eqn.27, and how this may differ from the strategies
 413 used in [11] to give numerical values to the terms in Eqn.2, but this is an area for future work.

414 Eqn.27 is our clarification of the algebraic structure of the ABC Score formula. Eqn.27 rigorously follows
 415 if enhancers collectively regulate a gene according to the IAC model (Assumptions 1 to 6).

416 Enhancer additivity and departures from it

417 The default model, satisfying Assumptions 1 to 4, exhibits *response additivity*, as shown by Eqn.19: the
 418 response of the gene to all the enhancers acting collectively is just the sum of the responses to each individual
 419 enhancer. When the default model also obeys Assumption 5 of deletion fidelity, then response additivity has
 420 a counterpart in the deletion effect, as defined in Eqn.20. This allows us to rigorously define the properties
 421 of *super-additivity* and *sub-additivity*. These departures from the properties of the default model may be
 422 helpful to interpret the effects of experimental perturbations, such as genetic deletions or CRISPRi, in which

423 subsets of enhancers are prevented from influencing a gene and the effect of these perturbations on the gene
424 expression response is measured.

425 With Assumptions 1 to 5, if U_1, \dots, U_m are pairwise disjoint subsets of enhancers, so that $U_i \subseteq$
426 $\{e_1, \dots, e_N\}$ and $U_i \cap U_j = \emptyset$ when $i \neq j$, then it follows from Eqn.24 that the effect of deleting all
427 the subsets together is just the sum of the individual deletion effects,

$$\Delta(G; U_1 \cup \dots \cup U_m) = \Delta(G; U_1) + \dots + \Delta(G; U_m). \quad (28)$$

428 We refer to this property as *deletion additivity*. Furthermore, it is evident from Eqn.24 that, if all the
429 enhancers are deleted, so that $U_1 \cup \dots \cup U_m = \{e_1, \dots, e_N\}$, then the *total deletion effect* must be 1,

$$\Delta(G; U_1) + \dots + \Delta(G; U_m) = \Delta(G; \{e_1, \dots, e_N\}) = 1. \quad (29)$$

430 Assuming deletion fidelity, the total deletion effect being 1 is equivalent to the response additivity in Eqn.19.
431 A special case of Eqn.29 arises if all enhancers are deleted individually, when, once again, the total deletion
432 effect is 1,

$$\Delta(G; e_1) + \dots + \Delta(G; e_N) = 1. \quad (30)$$

433 Now suppose that a gene g is regulated by N enhancers, e_1, \dots, e_N , each enhancer is modeled by the
434 graph G_i and the regulatory graph of g , G , has the product structure, $G_1 \times \dots \times G_N$. The labels in G need
435 not be related to those of the component graphs G_i , so that G need not be the product graph $G_1 \otimes \dots \otimes G_N$.
436 We can no longer calculate $R(G)$ in terms of $R(G_i)$. However, we can still define through Eqn.20 the
437 deletion effect $\Delta(G; U)$ for any collection $U \subseteq \{e_1, \dots, e_N\}$ of enhancers. We say that g exhibits response
438 *super-additivity* if,

$$R(G) > R(G_1) + \dots + R(G_N). \quad (31)$$

439 In terms of the deletion effect, this corresponds to when a collective deletion has less effect than the sum of
440 the individual deletions, so that,

$$\Delta(G; U_1 \cup \dots \cup U_m) < \Delta(G; U_1) + \dots + \Delta(G; U_m). \quad (32)$$

441 Similarly, g exhibits response *sub-additivity* if,

$$R(G) < R(G_1) + \dots + R(G_N). \quad (33)$$

442 and this corresponds to the collective deletion having more effect than the sum of the individual deletions,

$$\Delta(G; U_1 \cup \dots \cup U_m) > \Delta(G; U_1) + \dots + \Delta(G; U_m). \quad (34)$$

443 Experimentally, response additivity [15, 23, 38–43], super-additivity [15, 21, 39–44] and sub-additivity
444 [38, 40, 41] have all been observed. Because the super-additive and sub-additive findings cannot be accounted
445 for by the default model with deletion fidelity, we next consider some extensions of this model that show
446 how such effects could arise.

447 Mechanisms beyond the default model

448 In the following sections, we examine two departures from the default model and consider their impact on
449 whether enhancers act additively (Eqn.19), super-additively (Eqn.31) or sub-additively (Eqn.33). This will
450 also illustrate how our modeling framework can be used to reason about different biological mechanisms.

451 Non-additivity in mRNA production rates

452 In the default model, the summation of production rates in Assumption 3 is crucial for the property of
453 enhancer additivity in Eqn.19. The production rate is a convenient abstraction that aggregates over many
454 underlying molecular mechanisms, such as RNA Polymerase recruitment, pausing and elongation. It is

conceivable that, when multiple enhancers jointly influence transcription, the resulting rate is a more complex function than simple addition [38]. Here, we consider the effect of dropping Assumption 3.

Let us assume that we have two enhancers, e_1 and e_2 , which are described by the graphs $H_1 = K_{a,1} \otimes K_{c,1}$ and $H_2 = K_{a,2} \otimes K_{c,2}$, respectively, as specified in Assumption 6 in the coarse-grained version of our default model. The overall regulatory graph is given by $H_1 \otimes H_2$, so that e_1 and e_2 remain independent (Assumption 2). Note that $(K_{a,1} \otimes K_{c,1}) \otimes (K_{a,2} \otimes K_{c,2})$ has a product hierarchy and its vertices are therefore indexed by tuples of tuples of the form,

$$((a_1, c_1), (a_2, c_2)). \quad (35)$$

Here, a_i and c_i , for $i = 1, 2$, are coordinates for activation and communication, respectively, which take the values 0 and 1 in all cases. The graphs H_1 and H_2 have mRNA production rates, r_1 and r_2 , respectively, as specified in Eqn.15 and mRNA degradation rate δ .

We now consider a copy-number graph, G^\diamond , whose regulatory graph is given by $(K_{a,1} \otimes K_{c,1}) \otimes (K_{a,2} \otimes K_{c,2})$ but whose production rates do not obey Assumption 3. Note that we use the same symbol, G^\diamond , for the regulatory graph and the copy-number graph and rely on the context to clarify which is meant. There are many ways to assign production rates to the vertices of G^\diamond which do not obey Assumption 3; here we consider one of the simplest possible ways. We define 3 subsets of vertices of G^\diamond in terms of the coordinates in Eqn.35: $W := \{a_1 = 1, c_1 = 1, a_2 = 1, c_2 = 1\}$, $U := \{a_1 = 1, c_1 = 1\} \setminus W$ and $V := \{a_2 = 1, c_2 = 1\} \setminus W$. We assign the production rate of vertices in U to be r_1 , of vertices in V to be r_2 and of the vertex in W to be $(1 + \mu)(r_1 + r_2)$; all other vertices have production rate 0. We can summarise these assumptions in the following table, which gives the production rate for each of the 16 states in G^\diamond in the coordinate system described by Eqn.35.

state	rate	state	rate
$((0, 0), (0, 0))$	0	$((1, 0), (0, 0))$	0
$((0, 0), (0, 1))$	0	$((1, 0), (0, 1))$	0
$((0, 0), (1, 0))$	0	$((1, 0), (1, 0))$	0
$((0, 0), (1, 1))$	r_2	$((1, 0), (1, 1))$	r_2
$((0, 1), (0, 0))$	0	$((1, 1), (0, 0))$	r_1
$((0, 1), (0, 1))$	0	$((1, 1), (0, 1))$	r_1
$((0, 1), (1, 0))$	0	$((1, 1), (1, 0))$	r_1
$((0, 1), (1, 1))$	r_2	$((1, 1), (1, 1))$	$(1 + \mu)(r_1 + r_2)$

We further assume that $\mu \geq -1$ to ensure that the production rate of W does not become negative. If $\mu = 0$, then Assumption 3 holds for G^\diamond but not otherwise. We also assume that G^\diamond has degradation rate δ .

We now calculate $R(G^\diamond)$. Using the Sanchez-Kondev theorem in Eqn.14 we have,

$$R(G^\diamond) = \frac{r_1 u_U^*(G^\diamond) + r_2 u_V^*(G^\diamond) + (1 + \mu)(r_1 + r_2) u_W^*(G^\diamond)}{\delta}. \quad (36)$$

Expanding the term on the right hand side of Eqn.36 and rearranging terms results in,

$$R(G^\diamond) = \frac{r_1 [u_U^*(G^\diamond) + u_W^*(G^\diamond)] + r_2 [u_V^*(G^\diamond) + u_W^*(G^\diamond)] + \mu(r_1 + r_2) u_W^*(G^\diamond)}{\delta}. \quad (37)$$

Using the fact that the pairs of sets U and W , and, V and W are disjoint gives,

$$R(G^\diamond) = \frac{r_1 u_{U \cup W}^*(G^\diamond) + r_2 u_{V \cup W}^*(G^\diamond) + \mu(r_1 + r_2) u_W^*(G^\diamond)}{\delta}. \quad (38)$$

We now note that, by definition, $U \cup W = \{a_1 = 1, c_1 = 1\}$ and $V \cup W = \{a_2 = 1, c_2 = 1\}$. Given the independence assumption on G^\diamond , we can apply Eqn.11 and have that $u_{\{a_1=1, c_1=1\}}^*(G^\diamond) = \tilde{\alpha}_1 \tilde{\gamma}_1$, $u_{\{a_2=1, c_2=1\}}^*(G^\diamond) = \tilde{\alpha}_2 \tilde{\gamma}_2$ and $u_W^*(G^\diamond) = \tilde{\alpha}_1 \tilde{\gamma}_1 \tilde{\alpha}_2 \tilde{\gamma}_2$. Substituting into Eqn.38, we have,

$$R(G^\diamond) = \frac{r_1 \tilde{\alpha}_1 \tilde{\gamma}_1 + r_2 \tilde{\alpha}_2 \tilde{\gamma}_2 + \mu(r_1 + r_2) \tilde{\alpha}_1 \tilde{\gamma}_1 \tilde{\alpha}_2 \tilde{\gamma}_2}{\delta}. \quad (39)$$

483 Given that

$$R(H_1) + R(H_2) = \frac{r_1 \tilde{\alpha}_1 \tilde{\gamma}_1 + r_2 \tilde{\alpha}_2 \tilde{\gamma}_2}{\delta}, \quad (40)$$

484 we see that e_1 and e_2 act additively for $\mu = 0$ (Eqn.19), act super-additively for $\mu > 0$ (Eqn.31) and act
485 sub-additively for $\mu \in [-1, 0)$ (Eqn.33).

486 Non-independence in regulatory transitions between enhancers

487 So far all the graphs we have considered obey regulatory independence as defined by Assumption 2: the state
488 of one enhancer does not affect the transitions or the rates of any other enhancer. Let us examine this more
489 closely for the IAC model, with just two enhancers, e_1 and e_2 , described by graphs H_1 and H_2 , respectively,
490 as in the previous subsection. According to Assumption 6, $H_1 = K_{a,1} \otimes K_{c,1}$ and $H_2 = K_{a,2} \otimes K_{c,2}$ and
491 the overall regulatory system is therefore described by the graph, $G = H_1 \otimes H_2 = (K_{a,1} \otimes K_{c,1}) \otimes (K_{a,2} \otimes$
492 $K_{c,2})$. Using Eqn.12 and the coordinate system in Eqn.35, we have that $u_{\{a_1=1\}}^*(G) = u_{\{a_1=1\}}^*(H_1)$ and
493 $u_{\{a_2=1\}}^*(G) = u_{\{a_2=1\}}^*(H_2)$. That is, the probability of activation of an enhancer does not depend on the
494 presence of the other enhancer. However, if probability of activation is measured by H3K27ac ChIP-Seq, there
495 is evidence that perturbation of a single enhancer can result in altered H3K27ac signal at distal enhancers
496 [15, 22, 45, 46]. If such changes at the distal enhancer are not caused by the perturbation method itself,
497 so that the perturbation obeys the fidelity conditions in Assumption 5, then such experiments suggest that
498 there may be non-independence between enhancers at the level of activation.

499 In order to model non-independence between enhancers, we will consider a gene to be modeled by the
500 graph G^\sharp , where G^\sharp is the product between an activation graph A^\sharp and a communication graph C , so that
501 $G^\sharp = A^\sharp \otimes C$ (Fig.5). A^\sharp has the structure $K_{a,1} \times K_{a,2}$, in which $K_{a,1}$ and $K_{a,2}$ are both present as the
502 subgraphs,

$$(0, 0) \xrightleftharpoons[l_{a,1}]{k_{a,1}} (1, 0) \quad \text{and} \quad (0, 0) \xrightleftharpoons[l_{a,2}]{k_{a,2}} (0, 1),$$

503 respectively (Fig.5a). The remaining labels, on the edges $(1, 0) \rightleftharpoons (1, 1)$, which specify the rates of activation
504 and deactivation of e_2 when e_1 is activated, and on the edges $(0, 1) \rightleftharpoons (1, 1)$, which specify the rates of
505 activation and deactivation of e_1 when e_2 is activated, can be arbitrary. For simplicity, we assume that
506 $C = K_{c,1} \otimes K_{c,2}$, (Fig.5b), but note that non-independence in communication could be considered similarly.
507 G^\sharp has the same structure as $H_1 \times H_2$, and thus still models the activation and communication statuses of
508 e_1 and e_2 , but using the form $G^\sharp = A^\sharp \otimes C$ allows us to clarify the independence relationships in G^\sharp . Under
509 this reorganization, the vertex in Eqn.35 is now described in a new coordinate system as,

$$((a_1, a_2), (c_1, c_2)). \quad (41)$$

510 We assume that G^\sharp has the same mRNA production rates as for the IAC model. In terms of the vertex
511 subsets $W = \{a_1 = 1, c_1 = 1, a_2 = 1, c_2 = 1\}$, $U = \{a_1 = 1, c_1 = 1\} \setminus W$ and $V = \{a_2 = 1, c_2 = 1\} \setminus W$, the
512 vertices in U have production rate r_1 , those in V have production r_2 and those in W have production rate
513 $r_1 + r_2$; all other vertices have production rate 0. We can summarise this in the following table, in which
514 the states are described by the coordinate system in Eqn.41.

state	rate	state	rate
$((0, 0), (0, 0))$	0	$((1, 0), (0, 0))$	0
$((0, 0), (0, 1))$	0	$((1, 0), (0, 1))$	0
$((0, 0), (1, 0))$	0	$((1, 0), (1, 0))$	r_1
$((0, 0), (1, 1))$	0	$((1, 0), (1, 1))$	r_1
$((0, 1), (0, 0))$	0	$((1, 1), (0, 0))$	0
$((0, 1), (0, 1))$	r_2	$((1, 1), (0, 1))$	r_2
$((0, 1), (1, 0))$	0	$((1, 1), (1, 0))$	r_1
$((0, 1), (1, 1))$	r_2	$((1, 1), (1, 1))$	$r_1 + r_2$

515 As in the previous section, we can use the Sanchez-Kondev theorem in Eqn.14 to calculate,

$$R(G^\sharp) = \frac{r_1 u_U^*(G^\sharp) + r_2 u_V^*(G^\sharp) + (r_1 + r_2) u_W^*(G^\sharp)}{\delta} \quad (42)$$

$$= \frac{r_1 u_{U \cup W}^*(G^\sharp) + r_2 u_{V \cup W}^*(G^\sharp)}{\delta} \quad (43)$$

$$= \frac{r_1 u_{\{a_1=1, c_1=1\}}^*(G^\sharp) + r_2 u_{\{a_2=1, c_2=1\}}^*(G^\sharp)}{\delta}. \quad (44)$$

516 Given that $G^\sharp = A^\sharp \otimes C$, the probabilities of the sets $\{a_1 = 1, c_1 = 1\}$ and $\{a_2 = 1, c_2 = 1\}$ factor according
517 to Eqn.11. Continuing from Eqn.44 we have,

$$R(G^\sharp) = \frac{r_1 u_{\{a_1=1\}}^*(A^\sharp) u_{\{c_1=1\}}^*(C) + r_2 u_{\{a_2=1\}}^*(A^\sharp) u_{\{c_2=1\}}^*(C)}{\delta} \quad (45)$$

$$= \frac{r_1 u_{\{a_1=1\}}^*(A^\sharp) \tilde{\gamma}_1 + r_2 u_{\{a_2=1\}}^*(A^\sharp) \tilde{\gamma}_2}{\delta}. \quad (46)$$

518 Eqn.46 shows that the labels of A^\sharp do not directly appear in $R(G^\sharp)$; they only affect $R(G^\sharp)$ through the
519 enhancer activation probabilities $u_{\{a_1=1\}}^*(A^\sharp)$ and $u_{\{a_2=1\}}^*(A^\sharp)$. As in the previous section, we note that the
520 sum of the individual enhancer responses is,

$$R(H_1) + R(H_2) = \frac{r_1 \tilde{\alpha}_1 \tilde{\gamma}_1 + r_2 \tilde{\alpha}_2 \tilde{\gamma}_2}{\delta}. \quad (47)$$

521 Comparing Eqns.46 and 47, we see that whether the enhancers act additively (Eqn.19), sub-additively
522 (Eqn.33) or super-additively (Eqn.31) depends on the terms

$$\tilde{\alpha}_1^+ := \left[u_{\{a_1=1\}}^*(A^\sharp) - \tilde{\alpha}_1 \right] \quad \text{and} \quad \tilde{\alpha}_2^+ := \left[u_{\{a_2=1\}}^*(A^\sharp) - \tilde{\alpha}_2 \right]. \quad (48)$$

523 $\tilde{\alpha}_1^+$ represents the change in the probability of activation of enhancer 1 due to the presence of enhancer 2, and
524 $\tilde{\alpha}_2^+$ represents the change in the probability of activation of enhancer 2 due to the presence of enhancer 1. If
525 both $\tilde{\alpha}_1^+$ and $\tilde{\alpha}_2^+$ are positive, then e_1 and e_2 act super-additively; if they are both negative, then e_1 and e_2
526 act sub-additively. If $\tilde{\alpha}_1^+$ and $\tilde{\alpha}_2^+$ are of different signs, then the enhancers may act super-additively or sub-
527 additively depending on the relative magnitude of these terms compared to $\tilde{\gamma}_1, \tilde{\gamma}_2, r_1$ and r_2 . Experimental
528 data in which both $\tilde{\alpha}_1^+$ and $\tilde{\alpha}_2^+$ have been measured is limited. There are experimentally observed instances
529 in which both of these terms are positive [15, 45] but the precise form that the graph A^\sharp takes in these cases
530 is unknown. We are unaware of experiments that have observed $\tilde{\alpha}_1^+$ and $\tilde{\alpha}_2^+$ of differing signs. Whether
531 the experimentally observed non-additivity between enhancers can be explained by the non-independence
532 between enhancers as described in this section is a future area of research.

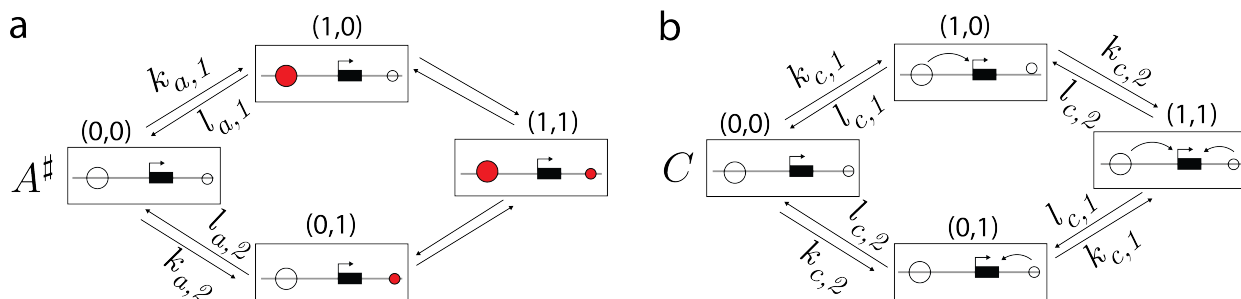


Figure 5: A model of non-independence in activation between enhancers. **(a)** The graph A^\sharp represents the activation components of each enhancer. A^\sharp has the structure of $K_{a,1} \times K_{a,2}$. Labels on the unmarked edges can be arbitrary. **(b)** The graph $C = K_{c,1} \otimes K_{c,2}$ represents the communication components of each enhancer.

533 Discussion

534 In this paper, we have introduced mathematical formulations of a *default model* and an *Independent-*
535 *Activation-Communication model* (IAC model) for how multiple enhancers collectively regulate a gene. The
536 default model encodes the notion that enhancers operate independently of each other (Assumptions 2 and 3).
537 At the same time, the default model imposes no assumptions on how the individual enhancers themselves
538 are working, at the level of transcription factors, co-regulators, chromatin, etc. They can be arbitrarily
539 complicated, so long as they operate within the Markovian setting that is commonly assumed for analysing
540 gene regulation. The default model assumptions imply that the collective response of a gene, as measured by
541 the mean mRNA level, is the sum of the responses coming from each enhancer individually, which we have
542 called *response additivity* (Eqn.19). The default model explains the mechanistic requirements for a gene to
543 exhibit this property and clarifies how ‘independence implies additivity’. We emphasize that independence
544 refers here to assumptions about gene regulatory mechanisms whereas additivity refers to the consequences
545 of those assumptions on steady-state gene expression. The default model supports the view that response
546 additivity is a reasonable baseline against which to assess the collective action of enhancers in regulating a
547 gene.

548 One of the advantages of the default model is that, because it is mathematically formulated, it allows
549 mechanistic departures from its assumptions to be systematically analysed. We have shown how depart-
550 ures from Assumptions 2 and 3 can give rise to response super-additivity (Eqn.32) as well as sub-additivity
551 (Eqn.34). As we have noted, response additivity [15, 23, 38–43], super-additivity [15, 21, 39–44] and sub-
552 additivity [38, 40, 41] have all been observed experimentally. The default model suggests the mechanistic
553 assumptions that could be experimentally tested to determine what underlies the observed response be-
554 haviour.

555 The IAC model is a special case of the default model that further assumes deletion fidelity, which allows
556 enhancers to be removed from the collective without influencing the remaining enhancers (Assumption 5),
557 and also assumes that individual enhancers can be described at a coarse-grained level in which they are
558 independently becoming activated and communicating their state to the gene (Assumption 6). Under As-
559 sumptions 1 to 6 for the IAC model, we derive a formula for the deletion effect of an individual enhancer
560 (Eqn.27) that shows a striking algebraic relationship to the ABC Score formula in Eqn.2. This relationship
561 suggests that the IAC model has accurately captured in mathematical terms the core intuitions behind the
562 ABC model from which the Score formula emerged [11].

563 A persistent conceptual theme that underlies the results reported here is that of independence. The
564 default model assumes that enhancers act independently, both in their regulatory state (Assumption 2) and
565 in their effect on mRNA production (Assumption 3). Furthermore the IAC model assumes that individual
566 enhancers become activated and communicating independently (Assumption 6). Our clarification of the
567 ABC Score formula thus arises from assuming independence *between* enhancers, along with independence
568 of activation and communication *within* each enhancer. The product construction on graphs and on graph
569 structures has been the key mathematical tool for rigorously defining independence, illustrating the value
570 of the graph-based linear framework for analysing gene regulation. We note that the concept of deletion
571 fidelity (Assumption 5) is also easily defined in the context of graphs.

572 Previous work used finite linear framework graphs to describe gene regulation [26]. Here, we have
573 introduced *copy-number graphs*, which have infinitely many vertices that keep track of both regulatory
574 states as well as the numbers of expressed mRNAs (Fig.2). Copy-number graphs allowed us to exploit
575 the Sanchez-Kondev theorem and calculate the mean mRNA number at steady state in terms solely of
576 the finite regulatory graph (Eqn.14). We therefore avoided dealing with infinite graphs despite relying on
577 them. Importantly, the linear framework also allows the unknown parameters within graphs to be treated
578 symbolically, so that conclusions may be drawn, as we saw above, without the need for assigning numerical
579 values to any of the parameters.

580 Beyond its utility in clarifying the ABC Score formula, the activation-communication coarse graining in
581 Assumption 6 provides an interesting lens through which to investigate enhancers. Many new experimental
582 technologies have emerged which allow perturbing entire enhancer sequences as a whole (as opposed to small
583 changes to DNA within an enhancer sequence). Such technologies include synthesizing and integrating long

584 DNA sequences [43, 44], modulating the genomic position of an enhancer [47–51], high throughput enhancer
585 perturbations with CRISPRi [11–13] and combining CRISPRi screens with rapid protein degradation [52].
586 By considering which perturbations affect, and do not affect, activation and communication, it may be
587 possible to probe the validity of the activation-communication coarse graining itself.

588 As noted in the Introduction, the ABC Score formula has been widely adopted for predicting enhancer-
589 gene connections. It has also been suggested that it could be combined with other predictive methods [53]
590 and that the ABC model could be used as a guiding principle in formulating other quantitative models
591 [54]. We believe the mathematical formulations that we have introduced here provide a foundation for such
592 efforts.

593 The ABC Score formula is quantitative (Eq.2) but the ABC model that gave rise to it is not a formal
594 mathematical model but, rather, an informal statement about the features, of enhancer activation and
595 contact, that are believed to be important in determining the response of a gene. Such informal models
596 play a critical role in biology but have the disadvantage that the underlying mechanistic requirements are
597 not clear. It is therefore difficult to know when the model can be applied and what can be deduced from it
598 when it does apply. In contrast, the mechanistic assumptions underlying our formal mathematical models
599 are precisely stated—Assumptions 1 to 4 for the default model and Assumptions 1 to 6 for the IAC model—
600 making it clear when the model can be applied and suggesting experimental tests to check the assumptions.
601 Moreover, if those assumptions are met, then the conclusions we have drawn, such as the response additivity
602 of the default model (Eqn.19) and the enhancer deletion formula for the IAC model (Eqn.27), are guaranteed
603 to hold as a matter of mathematical logic [55]. If those conclusions are not found experimentally, for example,
604 if response additivity is not found, then we know, as a matter of logic, that at least one of the assumptions
605 underlying the corresponding model does not hold. This understanding can, in turn, inform experiments to
606 determine where the departures from the assumptions occur. Such an approach allows a level of rigorous
607 reasoning about enhancer behaviour in gene regulation that is significantly harder to undertake with only
608 an informal quantitative model.

609 Mathematical theory has typically been introduced to analyse data, but the conceptual issues underlying
610 gene regulation are sufficiently intricate that theory may be necessary to understand the kinds of experiments
611 that are needed and how the data from them should best be interpreted [56]. Studies of the simple repression
612 motif in bacterial gene regulation may have already reached that point [57–60], as reviewed in [61]. The
613 foundation provided here, based on the linear framework, may offer similar opportunities in the eukaryotic
614 context. We believe our rigorous mathematical approach can play a significant role in investigating the
615 intricate interplay of enhancers in regulating gene expression.

616 Methods

617 A graph theory interpretation of the Sanchez and Kondev theorem

618 In this section we provide a proof of Eqn.14. We follow the Sanchez and Kondev approach described in
 619 [34] but present it using the graph theory notation and language used in this paper. Sanchez and Kondev
 620 provide in [34] a recurrence relation for all the moments of the mRNA probability distribution. A graph
 621 theory interpretation of these results, together with generalisations, will be presented in a separate paper;
 622 here we focus on the first moment only. We use bold face to denote matrices and vectors.

623 The steady-state distribution of a finite graph

624 Let G be a finite regulatory graph on the vertex set $V(G) = \{1, \dots, N\}$. We define the *Laplacian matrix* of
 625 G , $\mathcal{L} = \mathcal{L}(G)$, to be the $N \times N$ matrix,

$$\mathcal{L}(G)_{i,j} := \begin{cases} 0 & \text{if } j \neq i \text{ and } j \not\rightarrow i \\ \ell(j \rightarrow i) & \text{if } j \neq i \text{ and } j \rightarrow i \\ -\sum_{k \mid i \rightarrow k} \ell(i \rightarrow k) & \text{if } i = j. \end{cases} \quad (49)$$

626 As mentioned in the main text, G is equivalent to a continuous-time Markov process on the state space
 627 $\{1, \dots, N\}$ [25, 28, 32]. Let $u_i(t)$ be the probability that the process occupies state i at time t . Then the
 628 time evolution of the probability vector,

$$\mathbf{u}(t) := (u_1(t), \dots, u_N(t))^T,$$

629 is given by the master equation

$$\frac{d\mathbf{u}(t)}{dt} = \mathcal{L}(G)\mathbf{u}(t). \quad (50)$$

630 If G is strongly connected, then the kernel of $\mathcal{L}(G)$ is one dimensional, so there is a unique vector, $\mathbf{u}^*(G)$,
 631 such that $\mathcal{L}(G)\mathbf{u}^*(G) = 0$ and $u_1(G) + \dots + u_N(G) = 1$. $\mathbf{u}^*(G)$ is the steady-state probability distribution
 632 on G .

633 The master equation for a copy-number graph

634 Let $G \times P$ denote a copy-number graph with regulatory graph G , production rate vector $\mathbf{r} \in \mathbb{R}^N$ and
 635 degradation rate δ . Let $\mathbf{\Pi}$ be the diagonal matrix of production rates, $\mathbf{\Pi}_{i,i} = r_i$ and $\mathbf{\Pi}_{i,j} = 0$ when $i \neq j$,
 636 and let \mathbf{I} be the $N \times N$ identity matrix. Let

$$\mathbf{u}(p, t) := (u_{(1,p)}(G \times P; t), \dots, u_{(N,p)}(G \times P; t))^T,$$

637 be the vector of probabilities over the regulatory states with mRNA copy number p . It follows from the
 638 definition of the copy-number graph in the main text that $\mathbf{u}(p, t)$ satisfies the master equation,

$$\frac{d}{dt}\mathbf{u}(p, t) = \mathbf{\Pi}[\mathbf{u}(p-1, t) - \mathbf{u}(p, t)] + \delta\mathbf{I}[(p+1)\mathbf{u}(p+1, t) - p\mathbf{u}(p, t)] + \mathcal{L}(G)\mathbf{u}(p, t), \quad (51)$$

639 in which terms with arguments of $p-1$ are appropriately omitted when $p=0$. The first term of Eqn.51 arises
 640 from mRNA production, the second term from mRNA degradation and the third term from transitions in
 641 the regulatory graph. Eqn.51 can be rewritten as,

$$\frac{d}{dt}\mathbf{u}(p, t) = \mathbf{\Pi}\mathbf{u}(p-1, t) + \delta(p+1)\mathbf{u}(p+1, t) - [\mathbf{\Pi} + p\delta\mathbf{I} - \mathcal{L}(G)]\mathbf{u}(p, t). \quad (52)$$

642 We now let

$$\mathbf{q}(t) := \sum_{p=0}^{\infty} \mathbf{u}(p, t)$$

643 be the vector of marginal probabilities for the regulatory states. Proceeding from Eqn.52 we have,

$$\begin{aligned} \frac{d}{dt} \mathbf{q}(t) &= \sum_{p=0}^{\infty} \frac{d}{dt} \mathbf{u}(p, t) \\ &= \underbrace{\delta \mathbf{u}(1, t) - \mathbf{\Pi} \mathbf{u}(0, t) + \mathcal{L}(G) \mathbf{u}(0, t)}_{p=0} + \\ &\quad \underbrace{\mathbf{\Pi} \mathbf{u}(0, t) + 2\delta \mathbf{u}(2, t) - \mathbf{\Pi} \mathbf{u}(1, t) - \delta \mathbf{u}(1, t) + \mathcal{L}(G) \mathbf{u}(1, t)}_{p=1} + \\ &\quad \underbrace{\mathbf{\Pi} \mathbf{u}(1, t) + 3\delta \mathbf{u}(3, t) - \mathbf{\Pi} \mathbf{u}(2, t) - 2\delta \mathbf{u}(2, t) + \mathcal{L}(G) \mathbf{u}(2, t)}_{p=2} + \dots \end{aligned}$$

644 This is a telescoping sum which simplifies to

$$\frac{d}{dt} \mathbf{q}(t) = \mathcal{L}(G)(\mathbf{u}(0, t) + \mathbf{u}(1, t) + \mathbf{u}(2, t) + \dots) = \mathcal{L}(G) \mathbf{q}(t).$$

645 It follows that the steady-state marginal probability vector, \mathbf{q}^* , lies in the kernel of $\mathcal{L}(G)$ and must therefore
646 be equal to $\mathbf{u}^*(G)$,

$$\mathbf{q}^* = \mathbf{u}^*(G). \quad (53)$$

647 In other words, the steady-state marginal distribution of regulatory states in a copy-number graph is identical
648 to the steady-state distribution of regulatory states in a finite regulatory graph.

649 **Proof of Eqn.14**

650 We want to show that,

$$R(G \times P) = \sum_{(i,p)} p \cdot u_{(i,p)}^*(G \times P) = \frac{1}{\delta} \sum_{i \in V(G)} r_i \cdot u_i^*(G).$$

651 Let $\mathbf{u}^*(p)$ denote the steady-state probability distribution over the copy-number graph. (Note the distinction
652 with the marginal probability distribution over the regulatory states, $u^*(G) = \sum_p u^*(p)$.) Let

$$\boldsymbol{\mu}^* := \sum_{p=1}^{\infty} p \mathbf{u}^*(p)$$

653 be the corresponding steady-state average copy number vector. Evidently,

$$R(G \times P) = \mathbf{1}^T \boldsymbol{\mu}^*, \quad (54)$$

654 where $\mathbf{1}$ is the all-ones column vector of dimension N . Now let

$$\boldsymbol{\mu}(t) := \sum_{p=1}^{\infty} p \mathbf{u}(p, t)$$

655 be the time-dependent average copy-number vector. It follows from Eqn.52 that,

$$\frac{d}{dt} \boldsymbol{\mu}(t) = \sum_{p=1}^{\infty} p \frac{d}{dt} \mathbf{u}(p, t)$$

656

$$= \sum_{p=1}^{\infty} p \left(\mathbf{\Pi} \mathbf{u}(p-1, t) + (p+1) \delta \mathbf{u}(p+1, t) - (\mathbf{\Pi} + p\delta \mathbf{I} - \mathcal{L}(G)) \mathbf{u}(p, t) \right)$$

$$= \sum_{p=1}^{\infty} p \mathbf{\Pi} (\mathbf{u}(p-1, t) - \mathbf{u}(p, t)) + \sum_{p=1}^{\infty} p ((p+1) \delta \mathbf{u}(p+1, t) - p \delta \mathbf{u}(p, t)) + \sum_{p=1}^{\infty} p \mathcal{L}(G) \mathbf{u}(p, t) \quad (55)$$

The first summand in Eqn.55 can be simplified to,

$$\mathbf{\Pi} ((\mathbf{u}(0, t) - \mathbf{u}(1, t)) + 2(\mathbf{u}(1, t) - \mathbf{u}(2, t)) + \dots) = \mathbf{\Pi} (\mathbf{u}(0, t) + \mathbf{u}(1, t) + \dots) = \mathbf{\Pi} \mathbf{q}(t).$$

The second summand can be simplified to,

$$\delta ((2\mathbf{u}(2, t) - \mathbf{u}(1, t)) + 2(3\mathbf{u}(3, t) - 2\mathbf{u}(2, t)) + \dots) = -\delta \sum_{p=1}^{\infty} p \mathbf{u}(p, t) = -\delta \boldsymbol{\mu}(t).$$

And the third summand is evidently just $\mathcal{L}(G) \boldsymbol{\mu}(t)$. Combining these three simplifications, we see that,

$$\frac{d}{dt} \boldsymbol{\mu}(t) = \mathbf{\Pi} \mathbf{q}(t) - \delta \boldsymbol{\mu}(t) + \mathcal{L}(G) \boldsymbol{\mu}(t). \quad (56)$$

At steady state this becomes,

$$\delta \boldsymbol{\mu}^* - \mathcal{L}(G) \boldsymbol{\mu}^* = \mathbf{\Pi} \mathbf{q}^*.$$

Multiplying both sides $\mathbf{1}^T$, and recalling that,

$$\mathbf{1}^T \mathcal{L}(G) = \mathbf{0}^T \quad \text{and} \quad \mathbf{1}^T \mathbf{\Pi} = \mathbf{r}^T,$$

we find that,

$$\mathbf{1}^T \boldsymbol{\mu}^* = \frac{\mathbf{r}^T \mathbf{q}^*}{\delta}. \quad (57)$$

Using Eqns.53 and 54, we see that Eqn.57 becomes,

$$R(G \times P) = \frac{\mathbf{r}^T \mathbf{u}^*(G)}{\delta} \quad (58)$$

$$= \frac{1}{\delta} \sum_{i \in V(G)} r_i \cdot u_i^*(G). \quad (59)$$

as required. This completes the proof of Eqn.14.

Proof of response summation in the default model

In this section we prove Eqn.19 which shows that, within the default model, the collective response of all the enhancers is the sum of their individual responses. That is, if $G = G_1 \otimes \dots \otimes G_N$, then

$$R(G) = R(G_1) + \dots + R(G_N). \quad (60)$$

We consider the case with only two enhancers, $N = 2$, from which the general case follows easily. Recall from the Sanchez and Kondev formula in Eqn.14 that

$$R(G) = \frac{1}{\delta} \sum_{(i_1, i_2)} r_{(i_1, i_2)}(G) \cdot u_{(i_1, i_2)}^*(G). \quad (61)$$

Assumption 3 on response summation tells us that $r_{(i_1, i_2)}(G) = r_{i_1}(G_1) + r_{i_2}(G_2)$ and Eqn.8 for the product graph tells us that $u_{(i_1, i_2)}^*(G) = u_{i_1}^*(G_1) \cdot u_{i_2}^*(G_2)$. Substituting these expressions into Eqn.61 gives the following formula for $\delta \cdot R(G)$,

$$\sum_{(i_1, i_2)} (r_{i_1}(G_1) + r_{i_2}(G_2)) \cdot u_{i_1}^*(G_1) \cdot u_{i_2}^*(G_2).$$

674 We can perform the summation over (i_1, i_2) in any order, for instance by first summing over i_2 and then
675 summing over i_1 . This gives,

$$\sum_{i_1} \left(\sum_{i_2} r_{i_1}(G_1) \cdot u_{i_1}^*(G_1) \cdot u_{i_2}^*(G_2) + \sum_{i_2} r_{i_2}(G_2) \cdot u_{i_1}^*(G_1) \cdot u_{i_2}^*(G_2) \right). \quad (62)$$

676 In the inner left-hand sum over i_2 , the terms indexed by i_1 are constant and may be extracted from that
677 sum to give

$$r_{i_1}(G_1) \cdot u_{i_1}^*(G_1) \cdot \left(\sum_{i_2} u_{i_2}^*(G_2) \right). \quad (63)$$

678 Total probability always sums to 1, so that $\sum_{i_2} u_{i_2}^*(G_2) = 1$, and Eqn.63 reduces to

$$r_{i_1}(G_1) \cdot u_{i_1}^*(G_1). \quad (64)$$

679 Similarly, the inner right-hand sum in Eqn.62 may be written as

$$u_{i_1}^*(G_1) \cdot \left(\sum_{i_2} r_{i_2}(G_2) \cdot u_{i_2}^*(G_2) \right). \quad (65)$$

680 We recognise from Eqn.14 that the sum in brackets is the response of graph G_2 , so that Eqn.65 becomes,

$$u_{i_1}^*(G_1) \cdot \delta \cdot R(G_2). \quad (66)$$

681 We can now substitute Eqns.64 and 66 back into Eqn.62 to get,

$$\sum_{i_1} r_{i_1}(G_1) \cdot u_{i_1}^*(G_1) + \sum_{i_1} u_{i_1}^*(G_1) \cdot \delta \cdot R(G_2). \quad (67)$$

682 We recognise from Eqn.14 that the left-hand sum is δ times the response of graph G_1 . In the right-hand
683 sum, we can extract the terms that do not depend on i_1 and use once again that the total probability is 1.
684 This allows us to rewrite Eqn.67 as,

$$\delta \cdot R(G_1) + \delta \cdot R(G_2),$$

685 from which we conclude that, indeed,

$$R(G) = R(G_1) + R(G_2),$$

686 as claimed. This completes the proof of Eqn.19.

687 Symbolic computations

688 We have provided mathematical proofs for all of our results. However, many of our results were originally
689 discovered by exploration using computer algebra systems. Specifically, the Sage [62] computer algebra
690 system, and the SymPy [63] and NetworkX [64] Python packages were crucial for the development of this
691 paper.

692 Acknowledgements

693 JN, K-MN and JG were funded in part by NIH award R01GM122928; JN was also funded the NSF-Simons
694 Center for Mathematical and Statistical Analysis of Biology at Harvard University. We thank Zeba Wunder-
695 lich, Rosa Martinez-Corral, Jané Kondev and members of the Gunawardena lab for discussions and comments
696 on the manuscript.

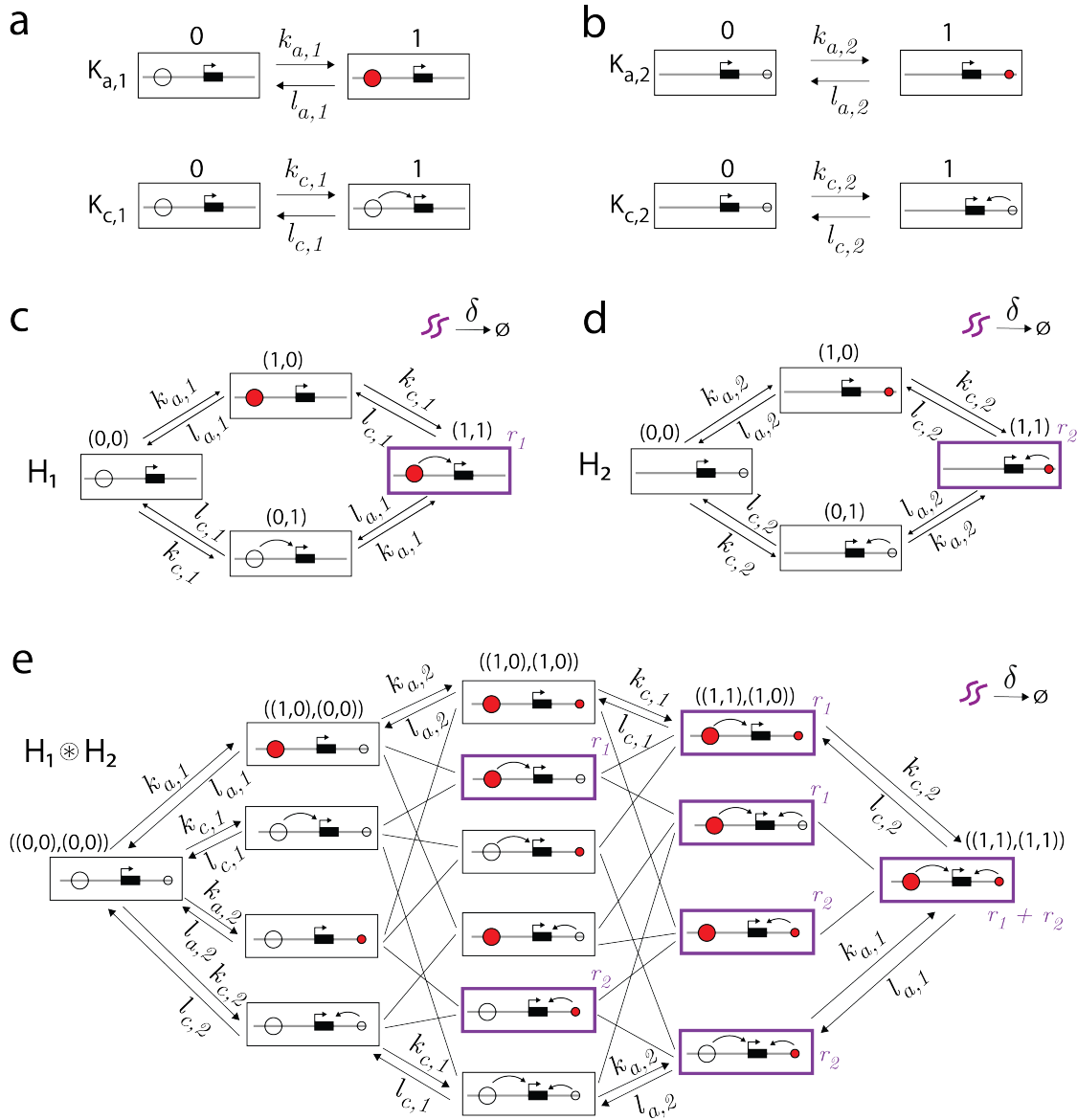


Figure S1: IAC model for $N = 2$ enhancers as a product graph of product graphs. **(a)** Graphs describing the activation and communication status of enhancer 1. **(b)** Graphs describing the activation and communication status of enhancer 2. **(c)** The graph H_1 whose underlying regulatory graph is $K_{a,1} \otimes K_{c,1}$. **(d)** The graph H_2 whose underlying regulatory graph is $K_{a,2} \otimes K_{c,2}$. **(e)** The graph $H_1 \otimes H_2$ satisfying the default model Assumptions 1-4 with components H_1 and H_2 . The regulatory graph of $H_1 \otimes H_2$ is given by $(K_{a,1} \otimes K_{c,1}) \otimes (K_{a,2} \otimes K_{c,2})$. Each vertex of this graph corresponds to the activation and communication statuses of both enhancers. For the vertices along the top of the graph, the product-graph binary notation is also provided using the coordinate system $((a_1, c_1), (a_2, c_2))$. Reverse edges and labels of most edges are omitted for clarity. Production states are highlighted in purple with corresponding production rates also in purple font.

697 References

- 698 1. Ptashne, M. & Gann, A. *Genes and Signals* (Cold Spring Harbor Laboratory Press, Cold Spring Harbor,
699 NY, USA, 2002).
- 700 2. Banerji, J., Rusconi, S. & Schaffner, W. Expression of a beta-globin gene is enhanced by remote SV40
701 DNA sequences. *Cell* **27**, 299–308 (1981).
- 702 3. Moreau, P., Hen, R., Wasylyk, B., Everett, R., Gaub, M. & Chambon, P. The SV40 72 base repair
703 repeat has a striking effect on gene expression both in SV40 and other chimeric recombinants. *Nucleic
704 Acids Research* **9**, 6047–6068 (1981).
- 705 4. Levine, M. Transcriptional enhancers in animal development and evolution. *Current biology : CB* **20**,
706 R754–R763 (2010).
- 707 5. Lempradl, A., Pospisilik, A. & Penninger, J. M. Exploring the emerging complexity in transcriptional
708 regulation of energy homeostasis. *Nat. Reviews. Genet.* **16**, 665–81 (2015).
- 709 6. Indjeian, V. B., Kingman, G. A., Jones, F. C., Guenther, C. A., Grimwood, J., Schmutz, J., Myers,
710 R. M. & Kingsley, D. M. Evolving new skeletal traits by cis-regulatory changes in bone morphogenetic
711 proteins. *Cell* **164**, 45–56 (2016).
- 712 7. Corradin, O. & Scacheri, P. C. Enhancer variants: evaluating functions in common disease. *Genome
713 Medicine* **6**, 85 (2014).
- 714 8. Noonan, J. P. & McCallion, A. S. Genomics of long-range regulatory elements. *Annual Review of
715 Genomics and Human Genetics* **11**, 1–23 (2010).
- 716 9. Zhang, X. & Meyerson, M. Illuminating the noncoding genome in cancer. *Nature Cancer* **1**, 864–872
717 (2020).
- 718 10. Fulco, C. P., Munschauer, M., Anyoha, R., Munson, G., Grossman, S. R., Perez, E. M., Kane, M., Cleary,
719 B., Lander, E. S. & Engreitz, J. M. Systematic mapping of functional enhancer–promoter connections
720 with CRISPR interference. *Science* **354**, 769–773 (2016).
- 721 11. Fulco, C. P. *et al.* Activity-by-contact model of enhancer–promoter regulation from thousands of
722 CRISPR perturbations. *Nature Genetics* **51**, 1664–1669 (2019).
- 723 12. Schraivogel, D., Gschwind, A. R., Milbank, J. H., Leonce, D. R., Jakob, P., Mathur, L., Korbel, J. O.,
724 Merten, C. A., Velten, L. & Steinmetz, L. M. Targeted perturb-seq enables genome-scale genetic screens
725 in single cells. *Nature Methods* **17**, 629–635 (2020).
- 726 13. Reilly, S. K. *et al.* Direct characterization of cis-regulatory elements and functional dissection of complex
727 genetic associations using HCR–FlowFISH. *Nature Genetics* **53**, 1166–1176 (2021).
- 728 14. Nasser, J. *et al.* Genome-wide enhancer maps link risk variants to disease genes. *Nature* **593**, 238–243
729 (2021).
- 730 15. Gschwind, A. R. *et al.* An encyclopedia of enhancer-gene regulatory interactions in the human genome.
731 *bioRxiv* (2023).
- 732 16. Bick, A. G. *et al.* Inherited causes of clonal haematopoiesis in 97,691 whole genomes. *Nature* **586**,
733 763–768 (2020).
- 734 17. Xu, Z. *et al.* Structural variants drive context-dependent oncogene activation in cancer. *Nature* **612**,
735 564–572 (2022).
- 736 18. Bendl, J. *et al.* The three-dimensional landscape of cortical chromatin accessibility in Alzheimer’s
737 disease. *Nature Neuroscience* **25**, 1366–1378 (2022).
- 738 19. Chen, Q., Dai, J. & Bian, Q. Integration of 3D genome topology and local chromatin features uncovers
739 enhancers underlying craniofacial-specific cartilage defects. *Science Advances* **8**, eabo3648 (2022).
- 740 20. Naqvi, S. *et al.* Precise modulation of transcription factor levels identifies features underlying dosage
741 sensitivity. *Nature Genetics*, 1–11 (2023).

- 742 21. Shin, H. Y., Willi, M., Yoo, K. H., Zeng, X., Wang, C., Metser, G. & Hennighausen, L. Hierarchy within
743 the mammary STAT5-driven Wap super-enhancer. *Nature Genetics* **48**, 904–911 (2016).
- 744 22. Li, K. *et al.* Interrogation of enhancer function by enhancer-targeting CRISPR epigenetic editing. *Nature*
745 *Communications* **11**, 485 (2020).
- 746 23. Hay, D. *et al.* Genetic dissection of the α -globin super-enhancer in vivo. *Nature Genetics* **48**, 895–903
747 (2016).
- 748 24. Gunawardena, J. A linear framework for time-scale separation in nonlinear biochemical systems. *PLOS*
749 *ONE* **7**, e36321 (2012).
- 750 25. Mirzaev, I. & Gunawardena, J. Laplacian dynamics on general graphs. *Bulletin of Mathematical Biology*
751 **75**, 2118–2149 (2013).
- 752 26. Ahsendorf, T., Wong, F., Eils, R. & Gunawardena, J. A framework for modelling gene regulation which
753 accommodates non-equilibrium mechanisms. *BMC Biology* **12**, 102 (2014).
- 754 27. Nam, K.-M., Martinez-Corral, R. & Gunawardena, J. The linear framework: using graph theory to
755 reveal the algebra and thermodynamics of biomolecular systems. *Interface Focus* **12**, 20220013 (2022).
- 756 28. Nam, K.-M. & Gunawardena, J. The linear framework II: using graph theory to analyse the transient
757 regime of markov processes. *Front. Cell Dev. Biol* **11**, 1233808 (2023).
- 758 29. Popay, T. M. & Dixon, J. R. Coming full circle: on the origin and evolution of the looping model for
759 enhancer–promoter communication. *Journal of Biological Chemistry* **298** (2022).
- 760 30. Karr, J. P., Ferrie, J. J., Tjian, R. & Darzacq, X. The transcription factor activity gradient (TAG)
761 model: contemplating a contact-independent mechanism for enhancer–promoter communication. *Genes*
762 *& Development* (2021).
- 763 31. Hnisz, D., Shrinivas, K., Young, R. A., Chakraborty, A. K. & Sharp, P. A. A phase separation model
764 for transcriptional control. *Cell* **169**, 13–23 (2017).
- 765 32. Norris, J. R. *Markov Chains* (Cambridge University Press, Cambridge, UK, 1997).
- 766 33. Peccoud, J. & Ycart, B. Markovian modeling of gene-product synthesis. *Theoretical Population Biology*
767 **48**, 222–234 (1995).
- 768 34. Sánchez, Á. & Kondev, J. Transcriptional control of noise in gene expression. *Proceedings of the National*
769 *Academy of Sciences* **105**, 5081–5086 (2008).
- 770 35. Dukler, N., Gulko, B., Huang, Y.-F. & Siepel, A. Is a super-enhancer greater than the sum of its parts?
771 *Nature Genetics* **49**, 2–3 (2017).
- 772 36. Van Arensbergen, J., FitzPatrick, V. D., de Haas, M., Pagie, L., Sluimer, J., Bussemaker, H. J. & van
773 Steensel, B. Genome-wide mapping of autonomous promoter activity in human cells. *Nature Biotech-*
774 *nology* **35**, 145–153 (2017).
- 775 37. Weingarten-Gabbay, S., Nir, R., Lubliner, S., Sharon, E., Kalma, Y., Weinberger, A. & Segal, E.
776 Systematic interrogation of human promoters. *Genome Research* **29**, 171–183 (2019).
- 777 38. Martinez-Ara, M., Comoglio, F. & van Steensel, B. Large-scale analysis of the integration of enhancer-
778 enhancer signals by promoters. *eLife* **12**, RP91994 (2024).
- 779 39. Loubiere, V., Almeida, B. P. d., Pagani, M. & Stark, A. Developmental and housekeeping transcriptional
780 programs display distinct modes of enhancer-enhancer cooperativity in drosophila. *bioRxiv* (2023).
- 781 40. Bothma, J. P., Garcia, H. G., Ng, S., Perry, M. W., Gregor, T. & Levine, M. Enhancer additivity and
782 non-additivity are determined by enhancer strength in the Drosophila embryo. *eLife* **4**, e07956 (2015).
- 783 41. Scholes, C., Biette, K. M., Harden, T. T. & DePace, A. H. Signal Integration by Shadow Enhancers and
784 Enhancer Duplications Varies across the Drosophila Embryo. *Cell Reports* **26**, 2407–2418.e5 (2019).
- 785 42. Lam, D. D. *et al.* Partially redundant enhancers cooperatively maintain mammalian pomc expression
786 above a critical functional threshold. *PLOS Genetics* **11**, e1004935 (2015).

- 787 43. Brosh, R. *et al.* Synthetic regulatory genomics uncovers enhancer context dependence at the Sox2 locus.
788 *bioRxiv* (2022).
- 789 44. Blayney, J. W. *et al.* Super-enhancers include classical enhancers and facilitators to fully activate gene
790 expression. *Cell* **186**, 5826–5839.e18 (2023).
- 791 45. Thomas, H. F. *et al.* Temporal dissection of an enhancer cluster reveals distinct temporal and functional
792 contributions of individual elements. *Molecular Cell* **81**, 969–982.e13 (2021).
- 793 46. Mansour, M. R. *et al.* An oncogenic super-enhancer formed through somatic mutation of a noncoding
794 intergenic element. *Science* **346**, 1373–1377 (2014).
- 795 47. Zuin, J. *et al.* Nonlinear control of transcription through enhancer–promoter interactions. *Nature*, 1–7
796 (2022).
- 797 48. Rinzema, N. J. *et al.* Building regulatory landscapes reveals that an enhancer can recruit cohesin to
798 create contact domains, engage CTCF sites and activate distant genes. *Nature Structural & Molecular*
799 *Biology* **29**, 563–574 (2022).
- 800 49. Thomas, H., Feng, S., Huber, M., Loubiere, V., Vanina, D., Pitasi, M., Stark, A. & Buecker, C. Enhancer
801 cooperativity can compensate for loss of activity over large genomic distances. *bioRxiv* (2023).
- 802 50. Jensen, C. L., Chen, L.-F., Swigut, T., Crocker, O. J., Yao, D., Bassik, M. C., Ferrell, J., Boettiger, A.
803 & Wysocka, J. Long range regulation of transcription scales with genomic distance in a gene specific
804 manner. *bioRxiv* (2024).
- 805 51. Brückner, D. B., Chen, H., Barinov, L., Zoller, B. & Gregor, T. Stochastic motion and transcriptional
806 dynamics of pairs of distal DNA loci on a compacted chromosome. *Science* **380**, 1357–1362 (2023).
- 807 52. Guckelberger, P. *et al.* Cohesin-mediated 3D contacts tune enhancer-promoter regulation. *bioRxiv*
808 (2024).
- 809 53. De Almeida, B. P., Reiter, F., Pagani, M. & Stark, A. DeepSTARR predicts enhancer activity from
810 DNA sequence and enables the de novo design of synthetic enhancers. *Nature Genetics* **54**, 613–624
811 (2022).
- 812 54. Kim, S. & Wysocka, J. Deciphering the multi-scale, quantitative cis-regulatory code. *Molecular Cell.*
813 *Reimagining the Central Dogma* **83**, 373–392 (2023).
- 814 55. Gunawardena, J. Models in biology: ‘accurate descriptions of our pathetic thinking’. *BMC Biol.* **12**, 29
815 (2014).
- 816 56. Phillips, R. Theory in Biology: Figure 1 or Figure 7? *Trends in Cell Biology* **25**, 723–729 (2015).
- 817 57. Garcia, H. G. & Phillips, R. Quantitative dissection of the simple repression input–output function.
818 *Proceedings of the National Academy of Sciences* **108**, 12173–12178 (2011).
- 819 58. Jones, D. L., Brewster, R. C. & Phillips, R. Promoter architecture dictates cell-to-cell variability in
820 gene expression. *Science* **346**, 1533–1536 (2014).
- 821 59. Brewster, R. C., Weinert, F. M., Garcia, H. G., Song, D., Rydenfelt, M. & Phillips, R. The transcription
822 factor titration effect dictates level of gene expression. *Cell* **156**, 1312–1323 (2014).
- 823 60. Razo-Mejia, M., Barnes, S. L., Belliveau, N. M., Chure, G., Einav, T., Lewis, M. & Phillips, R. Tuning
824 transcriptional regulation through signaling: a predictive theory of allosteric induction. *Cell Systems* **6**,
825 456–469.e10 (2018).
- 826 61. Phillips, R., Belliveau, N. M., Chure, G., Garcia, H. G., Razo-Mejia, M. & Scholes, C. Figure 1 theory
827 meets figure 2 experiments in the study of gene expression. *Annual Review of Biophysics* **48**, 121–163
828 (2019).
- 829 62. The Sage Developers. *SageMath, the Sage Mathematics Software System (Version 9.6)* (2022).
- 830 63. Meurer, A. *et al.* SymPy: symbolic computing in Python. *PeerJ Computer Science* **3**, e103 (2017).
- 831 64. Hagberg, A. A., Schult, D. A. & Swart, P. J. *Exploring network structure, dynamics, and function using*
832 *NetworkX* in *Proceedings of the 7th Python in Science Conference* (Pasadena, CA USA, 2008), 11–15.