

This is the accepted manuscript made available via CHORUS. The article has been published as:

Graph's Topology and Free Energy of a Spin Model on the Graph

Jeong-Mo Choi, Amy I. Gilson, and Eugene I. Shakhnovich

Phys. Rev. Lett. **118**, 088302 — Published 24 February 2017

DOI: [10.1103/PhysRevLett.118.088302](https://doi.org/10.1103/PhysRevLett.118.088302)

Graph's Topology and Free Energy of a Spin Model on the Graph

Jeong-Mo Choi,^{*} Amy I. Gilson,[†] and Eugene I. Shakhnovich[‡]
*Department of Chemistry and Chemical Biology, Harvard University,
 12 Oxford Street, Cambridge, Massachusetts 02138, USA*
 (Dated: February 2, 2017)

In this work we investigate a direct relationship between a graph's topology and the free energy of a spin system on the graph. We develop a method of separating topological and energetic contributions to the free energy, and find that considering the topology is sufficient to qualitatively compare the free energies of different graph systems at high temperature, even when the energetics are not fully known. This method was applied to the metal lattice system with defects, and we found that it partially explains why point defects are more stable than high-dimensional defects. Given the energetics, we can even quantitatively compare free energies of different graph structures via a closed form of linear graph contributions. The closed form is applied to predict the sequence space free energy of lattice proteins, which is a key factor determining the designability of a protein structure.

Over the last thirty years, graph theory has been applied to the study of various networks, including protein interaction networks, neural networks, and the World Wide Web [1–4]. While the interplay between network structure and dynamics has attracted considerable attention [5], the equilibrium characteristics of network have also been studied in various contexts. In this work, we will focus on a spin model on a graph, which has a wide range of applications from biochemical and immune network behaviors [6, 7] to social network phenomena [8], and study an analytical relationship between a system's graph topology and its free energy. Many of previous studies assume either a randomly generated graph [9], or a random Hamiltonian [10–13]. Instead of assuming a random variable, we assume a given (and arbitrary) interaction energy matrix over different spin states. This reduces generality on coupling, but by this we can investigate any type of graphs, whether dense or sparse.

Consider a simple graph [14] of N nodes. The graph connectivity is described by the adjacency matrix A , whose element A_{ij} is 1 when there is a link between nodes i and j , and $A_{ij} = 0$ otherwise. Each node is in one of M possible spin states. The Hamiltonian, \mathcal{H} , is defined as the summation of energetic contributions over all links, each of whose energy is determined by the states of its two terminal nodes. Note that orphan nodes do not contribute energetically, by definition. Now, the Hamiltonian can be written formally as

$$\mathcal{H} = \frac{1}{2} \sum_{i,j}^{N,N} A_{ij} E_{s(i)s(j)}, \quad (1)$$

where E is the energy matrix (whose element can be positive or negative) and $s(i)$ is the state of node i . The high-temperature expansion of the partition function $Z(\beta) = \sum_{\{s\}} e^{-\beta \mathcal{H}}$ over all possible state configurations is

$$Z(\beta) = \sum_{\{s\}} 1 - \beta \sum_{\{s\}} \mathcal{H} + \frac{\beta^2}{2!} \sum_{\{s\}} \mathcal{H}^2 - \dots \quad (2)$$

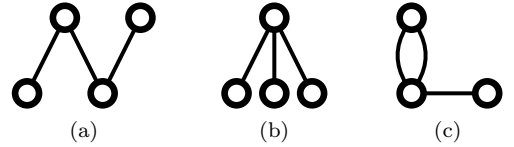


FIG. 1: Examples of 3-link multigraphs. a is a parent graph of c (*i.e.*, c is obtained by a node contraction operation on a), and b is also a parent graph of c, but there is no such relationship between a and b.

In the higher-order terms, there are summations of various products of A and E elements. To systematically visualize those terms, we introduce a multigraph g (where multiple links between any pair of nodes are allowed) with no orphan nodes (see Fig. 1 for examples), and we use its nodes as dummy variables of the summations. For each multigraph g , we define two quantities $[g]$ and $E(g)$ as follows:

$$[g] = \sum_{\text{nodes}} \prod_{k=1}^{n(g)} A_{l_k}, \quad (3)$$

$$E(g) = M^{-n(\text{nodes})} \sum_{\text{nodes}} \prod_{k=1}^{n(g)} E_{l_k}, \quad (4)$$

where l_k indicates the link of index k , $n(g)$ is the number of links in graph g , and $n(\text{nodes})$ is the number of nodes in graph g . For example, for the graph shown in Fig. 1c,

$$\left[\text{graph 1c} \right] = \sum_{ijk} A_{ij} A_{ij} A_{jk}, \quad (5)$$

$$E \left(\text{graph 1c} \right) = M^{-3} \sum_{mnp} E_{mn} E_{mn} E_{np}. \quad (6)$$

Note that since A is a binary matrix, $[g]$ for a graph g with multiple links is equal to $[g_0]$ where g_0 is a simple graph constructed from g by converting all multiple links

in g to single links. For example, equation 5 becomes

$$\left[\text{loop} \right] = \sum_{ijk} A_{ij}^2 A_{jk} = \sum_{ijk} A_{ij} A_{jk} = \left[\text{two nodes connected by two links} \right]. \quad (7)$$

Next, we introduce a special form of graph operation: node contraction, which merges two different nodes while preserving the number of links [15]. If a graph h can be obtained by any number of node contraction operations on another graph g , we will call h a child graph of g , and g a parent graph of h . For example, in Fig. 1, graph a is a parent graph of c, and graph b is also a parent graph of c, but there is no parent-child relationship between a and b. Note that $[g]$ represents the total number of unique subgraphs of type g and of its child types on the given graph.

Then it can be shown (see **Appendix A**) that

$$Z(\beta) = M^N \exp \left\{ \sum_{\text{connected } g} \frac{(-\beta/2)^{n(g)}}{n(g)!} H(g)[g] \right\}, \quad (8)$$

where the summation is over all possible connected subgraphs g . $H(g)$ is defined as

$$H(g) = K(g)E(g) + \sum_{g' \in \mathcal{P}(g)} (-1)^{m(g,g')} K(g,g') K(g') E(g'), \quad (9)$$

where $K(g)$ is the combinatoric factor to construct graph g from $n(g)$ links, $K(g,g')$ is the combinatoric factor to generate graph g from graph g' by node contraction, $m(g,g')$ is the number of contraction operations required to construct g from g' , and $\mathcal{P}(g)$ is the set containing all parent graphs of graph g . Finally, the free energy [16] is

$$F(\beta) = -Nk_B T \ln M + \sum_{\text{connected } g} \tilde{F}(g, \beta), \quad (10)$$

where

$$\tilde{F}(g, \beta) = -\frac{1}{\beta} \frac{(-\beta/2)^{n(g)}}{n(g)!} H(g)[g]. \quad (11)$$

Note that free energy depends only on global properties of a graph (such as $n(g)$ and $[g]$). Thus, all graphs with the same set of global properties share identical $\tilde{F}(g, \beta)$, independent of their local connectivity distributions.

One advantage of equations 10 and 11 is that the graph topology (which determines $[g]$) is now unlinked from detailed energetics (which determines $H(g)$). Hence, even without knowing the exact energy matrix E , it is possible to compare $[g]$ values from different structures and, in some cases, we can determine which structure provides lower free energy of the corresponding spin system. To illustrate this, let us consider two different graph systems, a chain graph of length N and a star graph with N leaves. They have the same numbers of nodes and links, but it

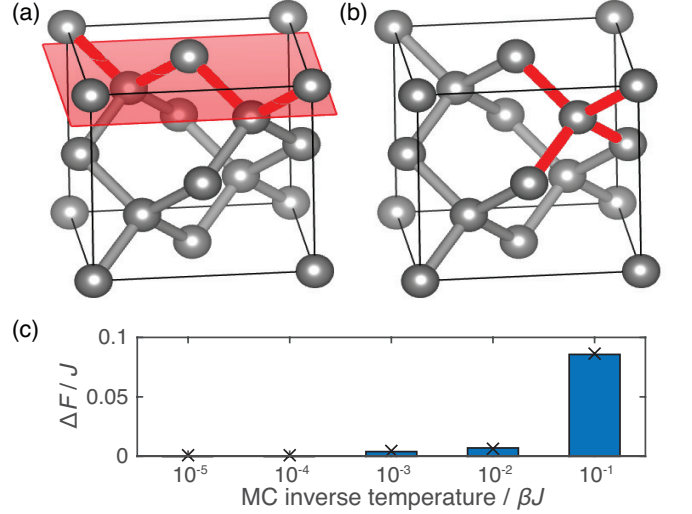


FIG. 2: (a-b) Two different types of defects, drawn by VESTA 3 [17]. Broken bonds are marked in red. (a) Diamond structure with a grain boundary indicated by a red plane. (b) Same structure with a vacancy defect. (c) Difference between MC free energies of two Si-Ge alloy lattice systems with different types of defects, as a function of inverse temperature. Here, the y axis shows $\Delta F = F(\text{vacancy defect}) - F(\text{grain boundary})$ for the entire system.

can be shown that $[g]^{\text{chain}} < [g]^{\text{star}}$ holds in general (see **Appendix D**), so at a temperature high enough that, for the star graph, the infinite sum in equation 10 does not diverge and if the sum is negative (stable free energy), we can conclude that the star-graph spin system has lower free energy than its counterpart on a chain graph. Note that this qualitative result is independent of the details of the energy matrix. As a specific example, an Ising chain system always has higher free energy than an Ising star at any $T > 0$ regardless of the details of energetics (see **Appendix D**).

A possible application of this general analysis is to an alloy lattice system with defects. Here we consider two types of defects: planar defects (*i.e.*, grain boundaries) and point defects. We constructed a 3-dimensional diamond-like lattice structure in $3 \times 3 \times 3$ unit cells with a periodic boundary condition (216 lattice points). The first system contains a grain boundary modeled by a discontinuity on the (001) plane (see Fig. 2a). To simulate vacancy defects, we constructed a second system by removing lattice points randomly (see Fig. 2b) until the number of broken bonds was equal to the number of bonds broken at the grain boundary in the first system (the number of remaining bonds = 324). Using a similar argument as above, it can be analytically shown that $[g]$ is generally greater for the point vacancy system than for the grain boundary system.

We carried out a Monte Carlo (MC) simulation to

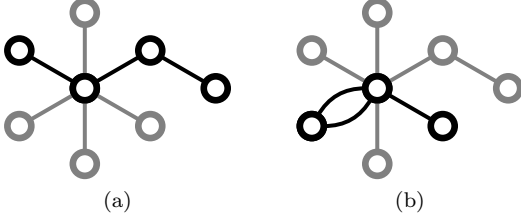


FIG. 3: Examples of subgraphs (black) contributing to the term $[g]$ where g is a 3-link path graph (Fig. 1a) on the given graph architecture (gray). (b) shows a child graph of a 3-link path graph (Fig. 1c).

check that the system with point vacancies indeed has lower free energy than that with a grain boundary. We considered lattice site occupation with the atom types $s(i) = \text{Si, Ge}$ as the possible states of lattice site i , and used the interatomic potential developed in previous works [18, 19]. Assuming equal bond lengths, $E(\text{Si, Si}) = -2.17J$, $E(\text{Ge, Ge}) = -1.93J$, $E(\text{Si, Ge}) = -2.04J$, where J is a constant. We tested five different temperatures, $10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}$, and 10^{-5} in units of βJ . We carried out 1,000 independent MC simulations for each temperature, and in each simulation we performed 1.1 million MC steps and neglected the first 0.1 million steps to achieve the system equilibrium (see Fig. S3 for a representative trajectory).

Fig. 2c shows the MC free energy difference between the two systems with different defect types as a function of βJ . In the high-temperature regime, the constant term and one-link term (proportional to the number of links) in equation 10 dominate so that the difference between the two systems is negligible. However, as inverse temperature βJ increases, the free energy difference between the two systems increases as well. The point defect system has lower free energy as expected, which is consistent with the experimentally known fact that point vacancies have thermal equilibrium concentrations whereas higher-dimensional defects cannot be formed spontaneously [20]. The entropic effect of multiple vacancy configurations and the stabilizing effect of structure relaxation have previously been used to explain this difference [21], but both factors were constant in our simulations so they cannot account for the free energy differences observed here. Also, the numbers of broken bonds are equal, meaning that the “surface areas” are the

same. Thus, this result implies that the stability of point defects (compared to line and planar defects) is partially due to the lattice topology itself.

We have hitherto considered qualitative differences between different graphs. Can we make a quantitative prediction of the system free energy? Although it is impossible to obtain a general closed form of equation 10, a closed form can be obtained for some special types of contributing graphs. We will focus on a path subgraph [22] (see Fig. 3a). Among tractable subgraphs, this type of graph has a significant contribution, because it does not have any connected parent graph and $[g] \geq [h]$ if g is a parent graph of h , so that the path graph provides one of the largest $[g]$ values among the connected graphs with the same number of links.

For a path subgraph g of length $n(g)$,

$$[g] = \sum A_{i_0 i_1} A_{i_1 i_2} \cdots A_{i_{n(g)-1} i_{n(g)}} = \text{su } A^{n(g)}, \quad (12)$$

where $\text{su } A$ is defined as the element sum of A , *i.e.* $\sum_{ij} A_{ij}$. Note that this term contains contributions from child subgraphs of g (see Fig. 3b). Similarly, we can express $E(g)$ for a path subgraph g by a relatively simple form,

$$E(g) = M^{-n(g)-1} \text{su } E^{n(g)}, \quad (13)$$

and we will define $\tilde{F}_{\text{path}}(\beta)$ as the summation of free energy contributors corresponding to path graphs:

$$\tilde{F}_{\text{path}}(\beta) = \sum_{\text{path } g} -\frac{1}{\beta} \frac{(-\beta/2)^{n(g)}}{n(g)!} K(g) E(g) [g]. \quad (14)$$

By using matrix diagonalization (see **Appendix E**), it can be shown that

$$\tilde{F}_{\text{path}}(\beta) = \frac{1}{2M^2} \sum_{i,j}^{N,M} \frac{|c_i|^2 |d_j|^2 \lambda_i \mu_j}{1 + \beta \lambda_i \mu_j / M}, \quad (15)$$

for $Mk_B T > \max(|\lambda_i|) \max(|\mu_j|)$. Here $\{\lambda_i\}$ and $\{\mu_j\}$ respectively represent the spectra of A and $\epsilon = E - \sum_{ij} E_{ij}/M^2$, and their corresponding eigenvector sets are respectively $\{|i\rangle_A\}$ and $\{|j\rangle_\epsilon\}$. We use inner products $c_i = \langle \mathbf{1} | i \rangle_A$ and $d_j = \langle \mathbf{1} | j \rangle_\epsilon$, by denoting an all-ones vector by $|\mathbf{1}\rangle$.

The free energy formula with the path graph factor is

$$F(\beta) = -Nk_B T \ln M + \frac{1}{2} E_0 \cdot \text{su } A - \frac{\beta}{4M^2} \left(\text{tr } \epsilon^2 - \frac{2}{M} \text{su } \epsilon^2 \right) \text{tr } A^2 + \frac{1}{2M^2} \sum_{i,j} \frac{|c_i|^2 |d_j|^2 \lambda_i \mu_j}{1 + \beta \lambda_i \mu_j / M} + \cdots, \quad (16)$$

where $\text{tr } A$ indicates the trace of A , while a simple linear approximation of equation 10 (see **Appendix F**) gives

$$F(\beta) = -Nk_B T \ln M + \frac{1}{2} E_0 \cdot \text{su } A - \frac{\beta}{4M^2} \left\{ \left(\text{tr } \epsilon^2 - \frac{2}{M} \text{su } \epsilon^2 \right) \text{tr } A^2 + \frac{2}{M} \text{su } \epsilon^2 \text{su } A^2 \right\} + \mathcal{O}(\beta^2). \quad (17)$$

To illustrate the utility of those approximate formulae, let us consider an example from biophysics. The designability of a protein structure is defined as the number of sequences that fold into the given structure as their lowest energy state. Biophysicists have been used this concept to investigate the principles of protein design and protein evolution (for review, see [23]). As previously discussed [24], there is a strong relationship between designability and sequence space free energy, *i.e.* the free energy of a heteropolymer in sequence space, instead of conformation space.

The Hamiltonian of a protein structure is given by

$$\mathcal{H} = \frac{1}{2} \sum_{i,j}^{N,N} A_{ij} E_{AA(i)AA(j)}. \quad (18)$$

Here A is called a contact matrix in the protein structure literature. Each element of A , A_{ij} , is 1 when residues i and j are nearest neighbors on the lattice but not adjacent on the protein backbone, and $A_{ij} = 0$ otherwise. E is an interaction matrix that contains interaction energies for every pair of amino acid types. N is the chain length, and $AA(k)$ is the amino acid type of residue k . This formula is analogous to the Hamiltonian for a spin model (equation 1; see [25]), so we can apply equation 16, or equation 17 to calculate the sequence space free energy at high sequence space temperature (where mutations can occur relatively frequently).

We studied the sequence space of a $3 \times 3 \times 3$ lattice protein structure (Fig. 4, upper left), whose intra-chain interaction network can be represented by a graph with 27 nodes and 28 links (the total number of non-covalent contacts). We used 10,000 representative structures among 103,346 maximally compact structures [26] to reduce the computational cost, following Heo *et al* [27]. We also employed two types of amino acids, and the interaction matrix represents hydrophobic and polar interactions:

$$E = \begin{bmatrix} -3J & J \\ J & 0 \end{bmatrix} \quad (19)$$

We designated an arbitrary chosen conformation as “native structure” and scanned all 2^{27} (about 1.3×10^8) possible sequences to compute $Z(\beta) = \sum_{\text{sequences}} e^{-\beta \mathcal{H}}$ and the corresponding $F(\beta)$ for a given β . We denote this F as the “exact sequence space free energy,” to distinguish it from a prediction made using equation 16 or 17, which we will call the “predicted sequence space free energy.” We repeated this calculation for all 10,000 structures.

Fig. 4 shows the sequence space free energy distributions over 10,000 lattice proteins at temperature $\beta J = 0.1$. The predicted values from equation 16 correlate strongly with the exact values (red, right axis), while predictions using the simpler approximation, equation 17, are not capable of discriminating between structures with different sequence space free energies (blue, left axis)

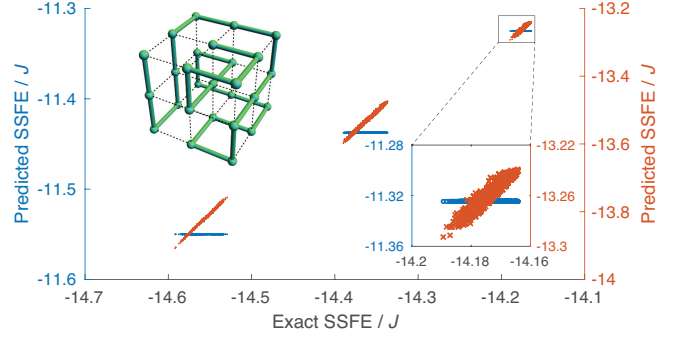


FIG. 4: Scatter plots of sequence space free energy (SSFE) distributions for 10,000 lattice protein structures, comparing exact and predicted SSFEs at $\beta J = 0.1$. Equation 16 (red) and equation 17 (blue) were used to calculate predicted SSFEs. (Upper left) a cartoon of a prototypical lattice protein. (Lower right) zoom of the boxed region.

[28]. It is because different structures share same $\text{su } A^2$ and $\text{tr } A^2$ values, and they can be distinguished only by considering higher-order terms. However, even in the former case, strict one-to-one correspondence does not hold between the exact and predicted values (lower right), because of contributions from higher-order terms that $\tilde{F}_{\text{path}}(\beta)$ does not capture. Note that the structures are mainly grouped by three different $\text{su } A^2$ values, implying that the system is still in the high-temperature regime, where higher-order terms do not dominate.

In this Letter, we presented an analytical method for calculating the free energy of a spin model on a simple graph. Through this approach, we find that the free energy contribution of the graph topology, realized by products of adjacency matrix elements, can be separated from energetic factors. The theory was illustrated by comparing chain and star graphs. Without specifying the interaction matrix, we showed that the star graphs are more stable than chain graphs in the high-temperature regime. The approach was then applied to lattice models of alloys with different defect types, which lead to different free energies, even when the systems had the same defect surface areas. We also showed that linear graphs are special in the sense that their infinite sum can be computed exactly, and this approach was applied to the protein design problem. The relative order of sequence space free energies of lattice proteins was predicted relatively accurately by the formula containing the infinite sum from the linear graph contribution, whereas a mere linear approximation could not discriminate among structures with the same $\text{su } A^2$ values. We hope that this theory will be expanded and applied to other graph-related problems in physics, from more complex spin systems to biological systems and also social networks.

We appreciate valuable comments from Erel Levine

and Kunok Chang. This work is supported by NIH grant GM068670.

* jeongmochoi@wustl.edu; Current address: Department of Biomedical Engineering, Washington University in St. Louis, 1 Brookings Drive, Saint Louis, Missouri 63130, USA

† amygilson@fas.harvard.edu

‡ shakhnovich@chemistry.harvard.edu

- [1] D. J. Watts and S. H. Strogatz, *Nature* **393**, 440 (1998).
- [2] R. Albert and A.-L. Barabási, *Rev. Mod. Phys.* **74**, 47 (2002).
- [3] R. Pastor-Satorras and A. Vespignani, *Evolution and Structure of the Internet: A Statistical Physics Approach* (Cambridge University Press, 2007).
- [4] S. Dorogovtsev and J. Mendes, *Evolution of Networks: From Biological Nets to the Internet and WWW* (OUP Oxford, 2013).
- [5] S. N. Dorogovtsev, A. V. Goltsev, and J. F. F. Mendes, *Rev. Mod. Phys.* **80**, 1275 (2008).
- [6] R. Könnig and R. Eils, *Bioinformatics* **20**, 1500 (2004).
- [7] E. Agliari, A. Annibale, A. Barra, A. C. C. Coolen, and D. Tantari, *Journal of Physics A: Mathematical and Theoretical* **46**, 415003 (2013).
- [8] C. Bisconti, A. Corallo, L. Fortunato, A. A. Gentile, A. Massafra, and P. Pellè, *Frontiers in Psychology* **6**, 1698 (2015).
- [9] A. Barrat and M. Weigt, *The European Physical Journal B - Condensed Matter and Complex Systems* **13**, 547 (2000).
- [10] B. Wemmenhove and A. C. C. Coolen, *Journal of Physics A: Mathematical and General* **36**, 9617 (2003).
- [11] E. Agliari and A. Barra, *EPL (Europhysics Letters)* **94**, 10002 (2011).
- [12] S. Franz, M. Leone, F. Ricci-Tersenghi, and R. Zecchina, *Phys. Rev. Lett.* **87**, 127209 (2001).
- [13] F. Guerra and F. L. Toninelli, *Journal of Statistical Physics* **115**, 531 (2004).
- [14] A simple graph is defined as a graph with no self-loops on any node and no multiple links between any pair of nodes. It may be either connected or disconnected.
- [15] S. Hartung and N. Talmon, in *Theory and Applications of Models of Computation: 12th Annual Conference, TAMC 2015, Singapore, May 18-20, 2015, Proceedings*, Lecture Notes in Computer Science, edited by R. Jain, S. Jain, and F. Stephan (Springer International Publishing, 2015) pp. 260–271.
- [16] **For tree graphs, free energy can be calculated iteratively. See Appendix C.**
- [17] K. Momma and F. Izumi, *J. Appl. Cryst.* **44**, 1272 (2011).
- [18] M. Laradji, D. P. Landau, and B. Dünweg, *Phys. Rev. B* **51**, 4894 (1995).
- [19] L. Cannavacciuolo and D. P. Landau, *Phys. Rev. B* **71**, 134104 (2005).
- [20] L. Priester, in *Grain Boundaries: From Theory to Engineering* (Springer Netherlands, Dordrecht, 2013) Chap. 5, pp. 135–146.
- [21] Y. Tateyama and T. Ohno, *Phys. Rev. B* **67**, 174105 (2003).
- [22] A path graph is a graph where two (terminal) nodes have vertex degree 1 and the other nodes have degree 2.
- [23] N. V. Dokholyan, in *Structural Bioinformatics*, edited by P. Bourne and J. Gu (John Wiley & Sons, 2009) Chap. 39, pp. 963–984.
- [24] J. L. England and E. I. Shakhnovich, *Phys. Rev. Lett.* **90**, 218101 (2003).
- [25] E. Shakhnovich and A. Gutin, *Protein Eng.* **6**, 793 (1993).
- [26] E. Shakhnovich and A. Gutin, *J. Chem. Phys.* **93**, 5967 (1990).
- [27] M. Heo, S. Maslov, and E. Shakhnovich, *Proc. Nat. Acad. Sci. U. S. A.* **108**, 4258 (2011).
- [28] **The Pearson correlation coefficients are 0.999 (eq. 16) and 0.998 (eq. 17) for the whole data set, but they reduce to 0.95 (eq. 16) and 0.0 (eq. 17) for the data in the zoomed box of Fig. 4.**