

# Combinatorial complexity and compositional drift in protein interaction networks

Eric J. Deeds<sup>1</sup>, Jean Krivine<sup>2</sup>, Jérôme Feret<sup>3</sup>, Vincent Danos<sup>4</sup>, Walter Fontana<sup>5,\*</sup>

**1** Center for Bioinformatics and Department of Molecular Biosciences, The University of Kansas, Lawrence KS 66047, USA

**2** Laboratoire PPS de l'Université Paris 7 and CNRS, F-75230 Paris Cedex 13, France

**3** Laboratoire d'Informatique de l'École normale supérieure, INRIA, ÉNS, and CNRS, 45 rue d'Ulm, F-75230 Paris Cedex 05, France

**4** School of Informatics, University of Edinburgh, Edinburgh, UK

**5** Department of Systems Biology, Harvard Medical School, 200 Longwood Avenue, Boston MA 02115, USA

\* E-mail: walter@hms.harvard.edu

## Abstract

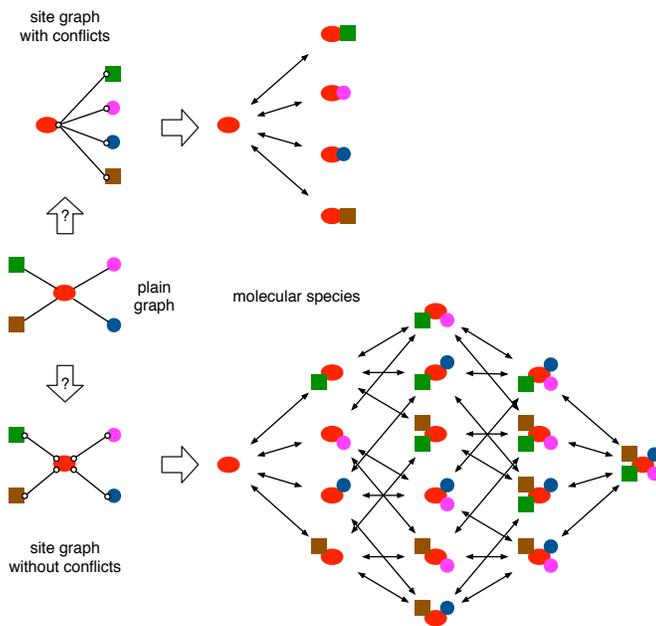
The assembly of molecular machines and transient signaling complexes does not typically occur under circumstances in which the appropriate proteins are isolated from all others present in the cell. Rather, assembly must proceed in the context of large-scale protein-protein interaction (PPI) networks that are characterized both by conflict and combinatorial complexity. Conflict refers to the fact that protein interfaces can often bind many different partners in a mutually exclusive way, while combinatorial complexity refers to the explosion in the number of distinct complexes that can be formed by a network of binding possibilities. Using computational models, we explore the consequences of these characteristics for the global dynamics of a PPI network based on highly curated yeast two-hybrid data. The limited molecular context represented in this data-type translates formally into an assumption of independent binding sites for each protein. The challenge of avoiding the explicit enumeration of the astronomically many possibilities for complex formation is met by a rule-based approach to kinetic modeling. Despite imposing global biophysical constraints, we find that initially identical simulations rapidly diverge in the space of molecular possibilities, eventually sampling disjoint sets of large complexes. We refer to this phenomenon as “compositional drift”. Since interaction data in PPI networks lack detailed information about geometric and biological constraints, our study does not represent a quantitative description of cellular dynamics. Rather, our work brings to light a fundamental problem (the control of compositional drift) that must be solved by mechanisms of assembly in the context of large networks. In cases where drift is not (or cannot be) completely controlled by the cell, this phenomenon could constitute a novel source of phenotypic heterogeneity in cell populations.

## Introduction

A large fraction of current data in molecular biology has been derived from the collation and curation of predominantly static types of data, such as genomic sequences and protein structures. However, at increasing rate, proteomic high-throughput methods, such as yeast two-hybrid assays, protein complementation assays, affinity purification with mass spectrometry, peptide phage display, and protein microarrays are yielding data about protein-protein interactions (PPI) whose significance resides in the system behavior they collectively generate [1–5]. In conjunction with more thorough biochemical measurements, these interaction data yield mechanistic statements ranging from less detailed, as in “*a phosphoepitope of EGFR binds strongly to the SH2/PTB domains of Grb2, Nck1, PI3K $\alpha$  and weakly to the SH2 domains of Grb10, Grb7, Nck2, Shp1*”, to more detailed, as in “*axin1 binds a region in the armadillo repeat of  $\beta$ -*

*catenin, if  $\beta$ -catenin is unphosphorylated at certain N-terminal residues.*” Unlike structural and genomic data types (“molecular nouns”), interaction fragments of this kind (“molecular verbs”) are fundamentally about process, and their broader meaning resides in the dynamic behavior of the large networks they generate.

High-throughput assays, such as yeast two-hybrid (Y2H), typically probe for pairwise binding between proteins in a highly impoverished context, lacking excluded volume and other effects that might influence interactions when the proteins tested are bound to multiple others [2,6]. Interaction data of this kind are often rendered as a large graph in which nodes represent proteins and edges correspond to pairwise binding interactions reported by the assay. These graphs have been shown to possess statistical properties, such as bow-tie structure [7,8], approximately scale-free degree distributions [9] and small-world characteristics [10]. Yet, unlike road networks, the edges in PPI networks do not represent persistent physical connections between nodes, but rather summarize interaction *possibilities* that must be realized through physical binding events. The cumulative effect of such events results in a distribution of protein complexes that ultimately determines cellular behavior. Significant properties of PPI networks may therefore become apparent only by studying the behavior they induce in a population of proteins, which requires the development and analysis of dynamic models.



**Figure 1. Binding surfaces and complex formation.** Center: The traditional plain graph representation of a PPI network represents the binding capabilities of a hub protein (red) through several incident edges. The diversity of molecular species generated by these potential interactions depends on the extent to which they compete for binding surfaces (white circles), to which we refer as “sites”. These conflicts are best represented as a “site graph”, derived from a domain-level resolution of protein-protein interactions. We depict two extreme cases. Top: All interaction partners compete for the same site. Bottom: All interactions occur at different sites and are mutually compatible. In the language we deploy to represent processes based on protein-protein interactions, a site denotes a distinct interaction capability. A comparison between the scenarios depicted at the top and the bottom illustrates how combinatorial complexity is affected by binding conflicts.

The first problem in constructing a dynamic model from raw PPI data is the lack of sufficient structural information. For instance, it is a priori unclear whether a “hub” protein with many interactions in the PPI network employs just one surface or many surfaces. As Figure 1 indicates, the set of complexes in which such a protein could participate depends on this information, since it allows the distinction between individual interactions that are mutually compatible and those that are mutually exclusive. The Structural Interaction Network (SIN) of yeast [11] is a dataset that provides this needed level of resolution.

It is often assumed that the various domains of a protein interact independently of one another; that is, the capacity of a protein’s domain  $A$  to bind its various partners is independent of the binding state of domain  $B$  on that same protein. While such an assumption represents an extreme case, so too does the assumption that domain  $A$  can bind only when domain  $B$  is unbound, or an assumption that posits strict allosteric correlations among binding partners. In the absence of systematic and readily accessible knowledge about steric and allosteric constraints in large-scale protein interaction networks, we consider the case of complete independence (subject to general biophysical constraints discussed below) as a useful “what-if” scenario against which to assess the significance of departures from independence.

The independence assumption creates a major challenge for making and running a model of a PPI network: the number of possible complexes (i.e. unique molecular species) that the network can generate increases exponentially as the network grows, reaching astronomical numbers for biologically reasonable networks [12, 13] (see also Figure 5 below). This situation necessitates an implicit representation of interactions as *local rules*, since models based on the explicit representation of all molecular possibilities, such as systems of differential equations, are entirely unfeasible. In recent years, we and others have developed appropriate tools for the representation and simulation of combinatorially complex systems of this kind [14–20].

In this contribution, we join two critical components—a suitable dataset and a modeling methodology—to simulate a large slice of the SIN network. By taking into account the inherent combinatorial complexity of the network, we extend pioneering calculations by Maslov and Ispolatov [21]. We consider neither post-translational modifications nor synthesis and degradation processes, as the available SIN data is exclusively about binding. Our simulated systems therefore reach thermodynamic equilibrium, although we shall see that this seemingly peaceful picture does not do justice to the microscopic dynamics. The main motivation for studying a highly abstracted and thus somewhat fictitious biochemical system is threefold. First, the image of a causally unconstrained network of possibilities, as conjured up by Y2H, has been taken seriously enough to attract extensive statistical investigation [22–25] of its structural properties. It seems warranted, therefore, to complement such studies with an eye on the dynamical properties implied by a similarly unconstrained interpretation of Y2H data. Second, the dynamic behavior of such a network serves as a null model to understand the need for and the consequences of curtailing independence through, for example, post-translational modification and allosteric interaction. In other words, studying the dynamics of the null model identifies a type of problem that specific causal constraints might have evolved to address, as we argue in the “Discussion” section. Third, the simulation of SIN dynamics represents a challenging test case illustrating a number of concepts underlying recent rule-based modeling methodologies [13–15, 17, 20] that are applicable to more general situations.

## Methods

### Interaction network data

As mentioned above, in order to provide a more structural picture of protein interaction networks, Kim *et al.* [11] combined raw interaction data from high-throughput experiments with data regarding domain-domain interactions in solved protein structures. This “Structural Interaction Network”—or SIN—associates a surface or domain of a protein with each interaction, converting the traditional flat

graph into a site graph or domain-level interaction network of the type shown in Figure 1. We obtained the original SIN directly from the authors. It consists of 1106 distinct proteins and 3826 specific pairwise interactions (edges).

Two proteins belong to the same graph component if there is a path of edges connecting them. The SIN has several such components. The largest (or “giant”) component consists of 454 proteins and 2572 interactions. The giant component contains 41% of the nodes in the graph, but includes 67% of its interactions. It therefore exhibits a significantly higher edge density (i.e. the fraction of possible edges present),  $\rho \approx 0.025$ , than the rest of the graph,  $\rho \approx 0.0059$ . The second-largest component in the SIN has only 21 proteins and most of the other components consist of only 2 proteins, representing isolated dimerizations. Current computational power precludes simulation of the dynamics of the entire SIN. Since the giant component contains a majority of the SIN interactions (and most of the interesting structure), we focussed on this part of the graph.

Data on subcellular localization and copy number were obtained from the “yeastgfp database” described in [26, 27]. This database contains information for about 75% of the proteins in the SIN. Using this data, we determined compartment-specific subgraphs of the SIN, consisting of only those proteins and their interactions that co-occur in the same compartment. These subgraphs exclude proteins that are found in a compartment but do not interact with any of the other proteins in that compartment, since such proteins could not participate in any kind of binding dynamics in our simulations. The cytoplasmic subgraph of the SIN consists of 349 proteins and 689 reactions. If we restrict ourselves to just the cytoplasmic subgraph of the giant component (which contains 78% of the interactions), we obtain a system with 167 proteins and 539 reactions, shown in Figure 2, which defines the network we simulated. We call this cytoplasmic subgraph of the giant component of the SIN the “cytoplasmic SIN” or cSIN for short.

Although homomeric interactions (i.e. a protein interacting with itself on some site) are certainly common, no such interactions have been characterized for this particular set of proteins: the *Saccharomyces* Genome Database (SGD, <http://www.yeastgenome.org>) lists no homomeric physical interactions for proteins in the cSIN.

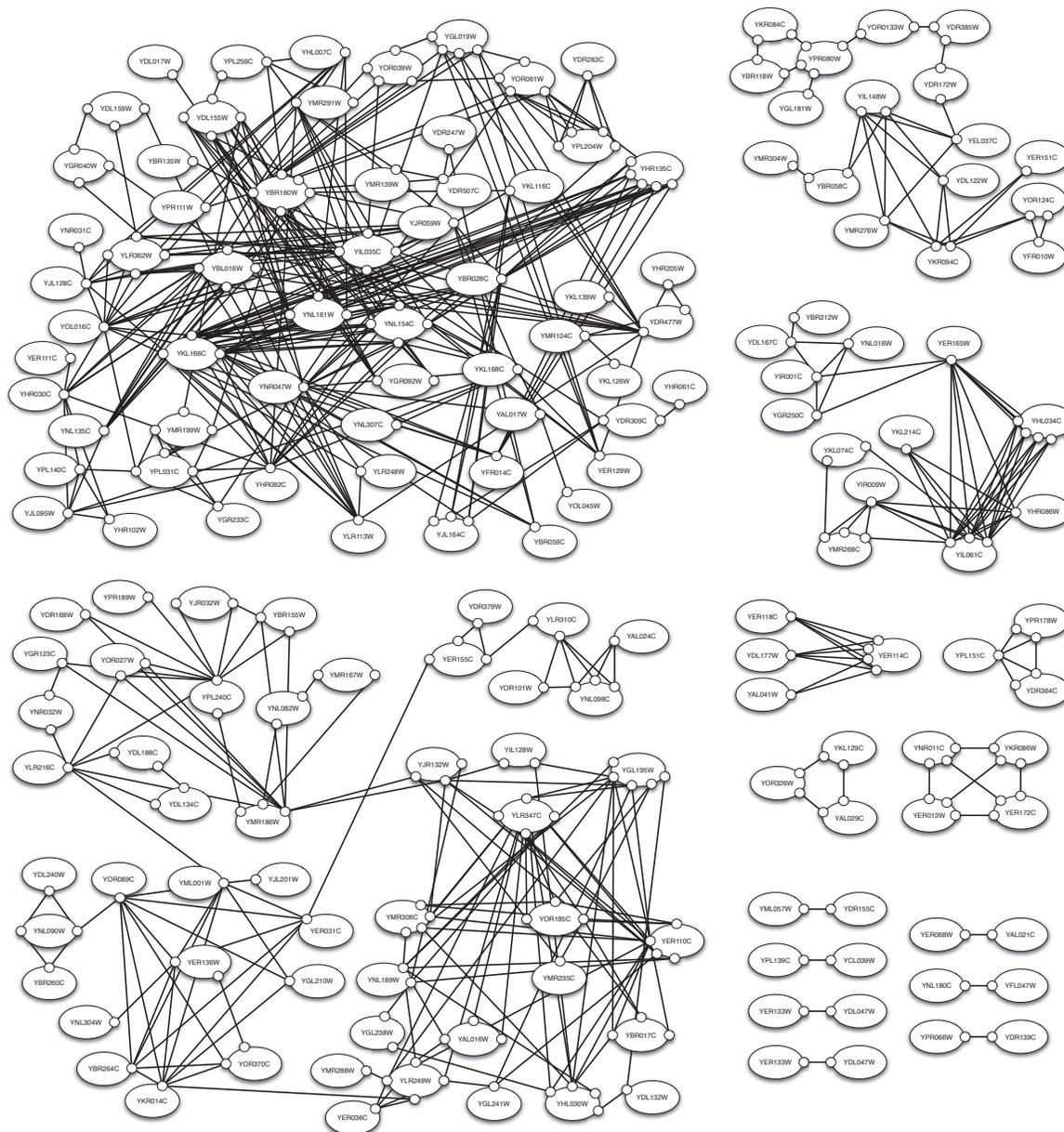
Copy numbers were assigned to each of these 167 proteins directly from the yeastgfp data [26]. In those cases where a protein is listed as existing in more than one compartment, assignment of a copy number to the cytoplasm becomes ambiguous. In the absence of data regarding the relative concentration of a given protein among compartments, we assumed that its concentration in each compartment is approximately equal. Since the cytoplasm represents the majority of the cell’s volume ( $\sim 85\%$  [28]), we simply assigned all copies of that protein to the cytoplasm. With this initial condition, the total number of individual protein agents present in each of our simulations was 2,908,889.

The localization and copy number data we used are based on measurements in asynchronous populations of cells [26, 27]. Our simulations do not take into account variations in copy number that might occur during the cell cycle [29–33]. However, only 13 of the 167 cSIN proteins exhibit strongly significant variations in expression level over the cell cycle, in the sense of being among the top 500 scoring yeast genes in a recent analysis [32]. Although changes in copy number during the cell cycle can clearly influence the types of complexes present in the cell [33], we leave consideration of these effects to future work.

A file with the complete set of interaction rules of the cSIN together with the initial condition is available as Supporting Information.

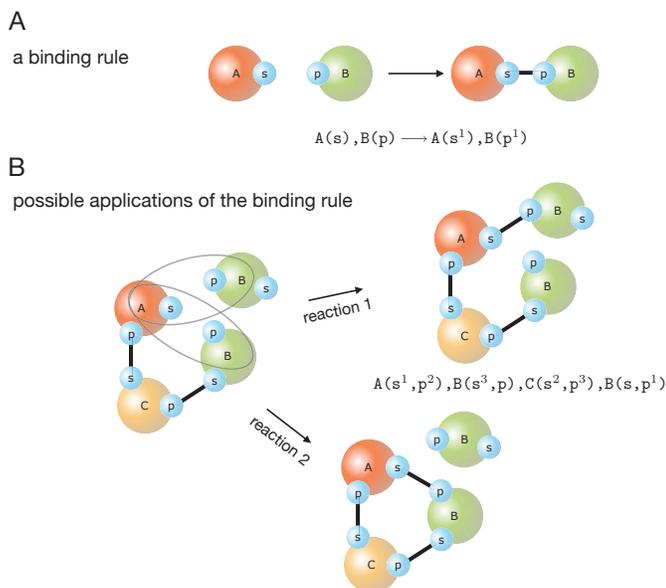
## Executable representation of the interaction network

A graph of *prima facie* independent binding interactions of the kind shown in Figure 2 permits a huge number of possible complexes (which we estimate in the “Results” section below). The vast number of possible molecular species rules out any modeling approach that requires their *a priori* enumeration. The only feasible simulation approach is one that replaces reactions between molecules with *local rules* that



**Figure 2. The network subject of this paper.** The graph of proteins, sites and interactions found in the cytoplasmic portion of the Structural Interaction Network (cSIN), as compiled by Kim et al [11]. The cSIN displays interactions at the level of domains or binding surfaces, making explicit which interactions compete for the same binding site. We refer to such a graph as a site graph. Its nodes are proteins (ovals), which are sets of sites (small circles on the ovals). Sites, rather than proteins, anchor the edges of this graph.

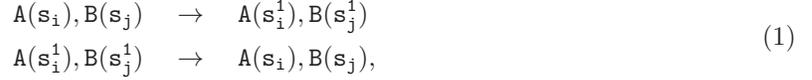
only specify which state modifications occur (in our case association or dissociation) and the sites on which these modifications depend (Figure 3). Reactions, on the other hand, must completely specify the binding state of each participating protein. A large set of reactions might express the same fundamental event in all of its possible contexts, whereas a rule can represent this entire family of reactions by specifying only the minimal context necessary for the event to occur. Rules can thus capture non-covalent association and dissociation of proteins or, more generally, post-translational modifications in a way that respects, as and when appropriate, the local quality of these interactions.



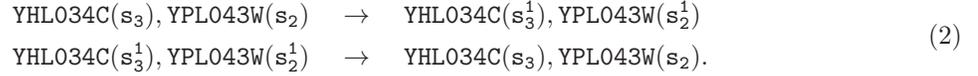
**Figure 3. Kappa rules.** **A:** A rule expresses a local mechanistic statement (of empirical or hypothetical origin) about a protein-protein interaction in terms of a rewrite directive plus a rate constant (not shown). The left hand side (LHS) of the rule consists of partially specified protein agents, and represents the contextual information necessary for identifying reaction instances that proceed according to the rule. The right hand side (RHS) expresses the actions that may occur when the conditions specified on the LHS are met in a reaction mixture. In this case, the rule specifies a binding action. Site graphs are represented in a simple syntax, explicated in Figure S1 of the Supporting Information. **B:** The rule in panel A can match the shown sample mixture of molecular species in two ways, giving rise to two possible reactions with different outcomes. Because of their local nature, Kappa-rules may apply in both a unimolecular and bimolecular situation. In general, such rules are given two rate constants (a first-order and a second-order constant), and the simulator will automatically generate the appropriate stochastic kinetics. However, in the present paper, global constraints prevent this ambiguity at the outset and the rules of the cSIN therefore necessitate only one rate constant (bimolecular for association and unimolecular for dissociation).

In representing and executing the cSIN, we follow our specification and implementation of a rule-based language, known as Kappa [14, 17, 18, 34–37], which is conceptually related to the Biological Network Generator Language (BNGL) [15, 16, 19, 20]; see section 1 of the Supporting Information. Rules that stipulate no other context than the domains involved in a binding or unbinding interaction between two proteins correspond exactly to the edges in the cSIN. We convert each edge into a pair of Kappa rules of

the kind



representing a binding (or unbinding) interaction between the  $i$ th site of protein A and the  $j$ th site of protein B. The superscript expresses a bond between the sites. For example:



Such rules of local interaction are then applied to a computational mixture consisting of a large graph whose nodes represent individual proteins and whose connected components represent protein complexes, much like the application of the rule in panel A of Figure 3 to the two-molecule mixture in panel B. Rule applications occur with probabilities in accordance with stochastic chemical kinetics, giving rise to a continuous-time Markov process implemented as detailed in [18, 19, 38] and summarized in the Supporting Information. At the start of a simulation, each protein is present with a number of copies derived from the previously mentioned empirical data, resulting in a total of  $\sim 3 \times 10^6$  individual protein agents.

## Affinities

In order to simulate the dynamics of a PPI network, we must assign to each (independent) binding reaction both an on-rate  $k_+$  (the rate constant for the first type of rule in equation 2) and an off-rate  $k_-$  (the rate constant for the second type of rule in equation 2). The dissociation constant,  $K_D \equiv k_-/k_+$ , is a measure of the strength or affinity of the corresponding interaction. Since high-throughput PPI experiments do not provide information about interaction strengths, we consider below three broad cases. The conversion into rate constants is discussed in the subsequent section.

### 1 — UNIFORM AFFINITIES.

Even when all of the binding reactions in the network have the same affinity, the question remains as to exactly *which* universal affinity to choose. The protein interaction strengths found in the PINT database exhibit an average affinity equivalent to a  $K_D$  of  $\sim 5$  nM [21, 39]. Since these interactions are obtained for a wide variety of proteins (many of which are not found in yeast and many of which represent mutated interaction pairs) and under a wide range of conditions (i.e. pH values and temperatures that are not necessarily characteristic of the yeast cytoplasm), it is difficult to interpret what this average value might mean for the cSIN. We therefore chose to look at a variety of  $K_D$  values: 10 nM, 100 nM and 1  $\mu$ M. The 10 nM case represents a set of fairly strong interactions (close to the average in PINT [21, 39]) and the 1  $\mu$ M case represents a set of fairly weak interactions.

### 2 — CONCENTRATION-BASED AFFINITIES (“EQUAL SATURATION”).

Even for strong interaction strengths (e.g. 10 nM), the log-normal distribution of protein concentrations observed within the cell causes reactions to operate at widely differing saturation levels. For instance, an interaction between two proteins at a concentration of  $\sim 1$   $\mu$ M will be highly saturated when assuming a  $K_D$  of 10 nM, while an interaction between two other proteins present at 0.1 nM will not be saturated at all. Following Maslov and Ispolatov [21], we consider a case in which each reaction in the network operates at approximately the same level of saturation. Consequently, we require the reaction affinities to vary with the (initial) reactant concentration as

$$K_D(i, j) = \frac{\max(C_i, C_j)}{20}, \quad (3)$$

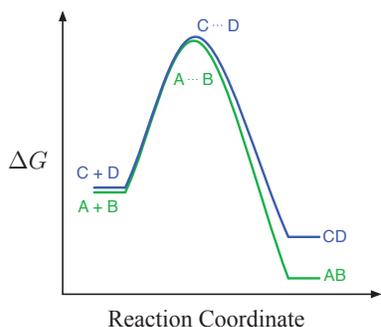
where  $K_D(i, j)$  is the dissociation constant of binding between proteins  $i$  and  $j$ , and  $C_x$  denotes the total concentration of protein  $x$  (obtained from experiment [40]). This method ensures that the overall binding saturation is essentially constant across reactions in the network when physiological concentrations are employed. The set of  $K_D$ 's obtained from equation 3 are log-normally distributed [40], and has recently been shown to represent a biologically and biophysically realistic case [41, 42].

### 3 — STRUCTURE-BASED AFFINITIES.

We can estimate binding affinities directly from the protein structures on which the interaction network is based [11]. Several studies have noted that the change in solvent-accessible, non-polar surface area that occurs on binding,  $\Delta\text{SASA}_{\text{NP}}$ , is linearly related to the free energy of association [43, 44]. To make use of this fact, we first re-constructed (as detailed in section 8.2 of the Supporting Information) the PPI network on the basis of the domain-domain interaction structures referenced in the most recent release of iPfam. We call this network the “cSIN2.” For each interaction in the cSIN2, we used the software package POPS [45] to determine the average  $\Delta\text{SASA}_{\text{NP}}$  taken over all the instances of that particular domain-domain interaction in iPfam. Using a recently published data set [44], we performed a linear regression to map  $\Delta\text{SASA}_{\text{NP}}$  into the corresponding free energy of binding  $\Delta G_b$ . Although the correlation in this case is certainly not perfect ( $R^2 = 0.47$ , see Figure S11 of the Supporting Information), the resulting equation provided us at least with a rough estimate of  $K_D$  (as  $\exp(\Delta G_b/RT)$ ) for each interaction in the cSIN2.

## Rate constants

We next describe the conversion of affinities into on- and off-rates. Let  $k_+(i, j)$  denote the rate constant of the binding reaction between proteins  $i$  and  $j$  (on-rate) and let  $k_-(i, j)$  denote the dissociation rate constant for that bond (off-rate). Since  $K_D(i, j) = k_-(i, j)/k_+(i, j)$  only constrains the ratio of the rates, we can choose either the on- or the off-rate arbitrarily and still satisfy a specified reaction affinity.



**Figure 4. Schematic free energy landscape.** The schematic shows the free energy landscape for a case in which differences in affinities are entirely represented by differences in off-rates. Here we have two different binding reactions: A binds B and C binds D. “A + B” and “C + D” represent the unbound states on the far left of the schematic reaction coordinate; the unbound states in this case have roughly the same free energy. The transition states (represented by “A ··· B” and “C ··· D”) also have approximately the same free energy; the change in free energy from the unbound state to the transition state is identical in both cases (giving identical values of  $k_+$ ). However, the bound states (“AB” and “CD”) exhibit very different free energies, and the difference in free energy change between the transition state and the bound state results in a much higher value of  $k_-$  for the C-D binding reaction compared to the A-B binding reaction.

In the present work, we constrain the on-rate to *always* have the same value, regardless of the  $K_D$ . When all reactions in the network have the same affinity, varying the global affinity (e.g. from 10 nM to 100 nM) thus amounts to varying the probability that bonds will be broken once they are formed. This means that the relative change in free energy between the unbound state and the binding transition state is the same for all reacting pairs; all that changes is the free energy of the bound state, as illustrated schematically in Figure 4. It appears reasonable [41, 42] that much of the differences in binding free energies across the network are due to differences in relative hydrophobicity. However, in cases where the transition state free energy includes significant electrostatic contributions, one might expect significant variance in both on- and off-rates [46].

Equipped with deterministic rate constants  $k$  for each of our reactions, we convert these into stochastic rate parameters  $\beta$ . A dimensional argument suggests that for a unimolecular unbinding reaction  $\beta_- = k_-$  in units of  $s^{-1}$ , while for a bimolecular binding reaction

$$\beta_+ = \frac{k_+}{N_A V}, \tag{4}$$

in units of  $\text{molecule}^{-1}\text{s}^{-1}$ , where  $k_+$  is the deterministic rate constant in units of  $\text{M}^{-1}\text{s}^{-1}$ ,  $N_A$  is Avogadro’s constant and  $V$  is the volume of the system in liters. Microscopically, the inverse volume dependence arises from converting the “collision volume” swept out by a moving molecule into a probability through division by the volume available to an encounter, i.e. the volume of the system [38]. A unimolecular reaction has no collision volume and therefore its stochastic rate is independent of the system volume.

Since the protein copy numbers used in our simulations were obtained for haploid yeast cells, we approximate the volume to be  $42 \mu\text{m}^3$ , or  $4.2 \times 10^{-14}$  L [47]. We set the on-rate  $\beta_+(i, j) = 0.01$  for all  $i, j$  in the network, which corresponds, by equation 4, to a deterministic on-rate of  $2.5 \times 10^8 \text{ M}^{-1}\text{s}^{-1}$ . Given the absence of empirical measurements, the value of  $k_+$  ( $\beta_+$ ) is not meant to be realistic. Interactions driven purely by hydrophobicity could have values  $\sim 10^6 - 10^7 \text{ M}^{-1}\text{s}^{-1}$  [48]. The time scales discussed in the “Results” section are estimated assuming this range of on-rates, but it is important to note that the actual on-rates observed in a living system might differ significantly. Hence, for our simulations, the unit of time is essentially arbitrary.

## Preventing polymerization

A local cSIN rule like equation 2 specifies the binding between specific domains of proteins A and B, without, however, specifying whether A and B are members of the same or distinct complexes. In the first case the interaction is intramolecular; in the second case it is intermolecular (Figure 3). When the underlying network site graph contains proper cycles (i.e. paths that start and end on the same protein node without touching a site twice), this ambiguity results in infinitely many possible rings and polymers. Without further constraints, mass action would lead to a prevalence of long polymers, but aside from cytoskeletal proteins (such as actin and tubulin) or prions there is no empirical information suggesting that proteins generally form non-covalent polymer chains. In our simulations we must, therefore, prevent or curb polymerization. We achieve this by employing *global* constraints, that is, constraints that are not expressed directly as executable rules, but as filters applied by the simulator at runtime. We implemented two scenarios that correspond to distinct structural interpretations of network cycles, which we summarize next. A detailed exposition can be found in sections 6 and 7 of the Supporting Information.

THE “STABLE RINGS” (SR) SCENARIO. We might imagine that the open chain  $R \equiv A - C - B$  (which, in the more precise notation of our formalism, reads  $A(\mathbf{s}, \mathbf{p}^1), C(\mathbf{s}^1, \mathbf{p}^2), B(\mathbf{s}^2, \mathbf{p})$ ) is structurally sufficiently constrained to readily form a cyclical complex by *intramolecular* binding between A and B. In this rationale, there is not enough physical room in  $R$  to accommodate another B in an *intermolecular* reaction with A. We refer to this scenario as “stable rings” (SR): In this case the binding site on A is assumed to be naturally occluded by the B already bound to C. In the SR scenario, ring-like structures are highly

stable [49] and form *immediately* whenever intramolecular ring closure is possible. A thermodynamic justification of this scenario is discussed in section 6.1 of the Supporting Information. Polymerization is thus prevented by the formation of stable rings and a constraint enforcing the excluded volume implied by the SR scenario (Figure S5 of the Supporting Information).

**THE “NO RINGS” (NR) SCENARIO.** Many steric constraints other than direct occlusion of A’s binding site for B might prevent the addition of a second B to R. We subsume these alternative geometries under the “no rings” (or NR) scenario. The NR scenario introduces a syntactical filter that simply prevents at runtime any form of polymerization by *fiat*, as detailed in Figure S6 and section 7.1 of the Supporting Information.

Neither the SR case nor the NR case is likely to represent the reality of complex formation in the cell. Some of the cycles in the contact map of the cSIN might represent SR complexes, others might follow the NR scenario or perhaps even give rise to polymers of limited size.

We assessed the validity of the cSIN and the soundness of our model by comparing our computational mixtures of complexes with Affinity Purification-Mass Spectrometry (AP-MS) experiments (section 9 of the Supporting Information). In discussing the computational results, we focus on the NR scenario since it provides slightly better overlap with experimental data.

## Results

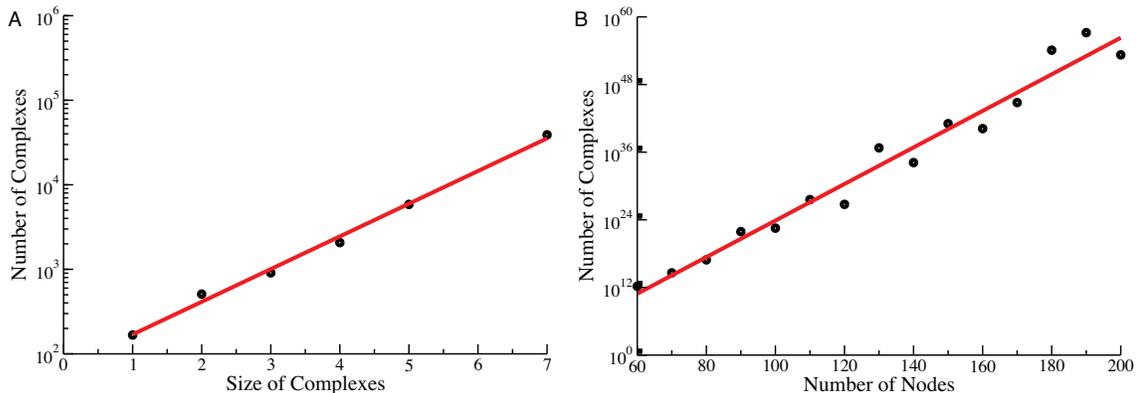
### Estimating the Number of Reachable Molecular Species

The number of distinct molecular species—the “reachable complexes” or “reachables” for short—that can, in principle, be generated with the interactions listed in the cSIN conveys a sense for the fraction of possibilities that a population of protein agents can access at any one time.

If an interaction network does not give rise to cyclical subgraphs, the set of reachables can be enumerated. If cycles are present, as is the case in the cSIN, the set of reachables, absent any constraints, is infinite due to polymerization. The cSIN contains many proper cycles (see Figure 2), which motivated the SR and NR scenarios described above. Since these constraints are not expressed as Kappa rules, but rather enforced at runtime, we were unable to compute the possibilities inherent in the cSIN other than by brute force enumeration stratified by complex size, as reported below. This strategy is feasible only up to a modest size. However, we can estimate the combinatorial complexity of the cSIN by constructing artificial *acyclic* interaction graphs with an edge density that matches the cSIN and for which we can count the number of complexes.

*Direct Enumeration by complex size.* The cSIN consists of 167 distinct proteins, and thus 167 unique monomers, and 539 dimers, since every interaction in the network can form a unique dimer. Starting from the set of dimers, we can create a set of trimers by taking a free site in every such dimer and adding a possible binding partner to form a trimer. Because of cycles in the contact map, such a procedure could easily produce multiple copies of the same complex; for instance, adding a C to the B of an A-B dimer produces the same A-B-C trimer as adding an A to the B of a B-C dimer. To avoid overcounting, we simply check for each new complex whether it has already been found and, if it has, we discard it. We prevent polymeric complexes by simply requiring that no agent type occurs twice in the same complex. This is a stricter criterion than the no-polymerization constraint of the NR scenario mentioned above. As such our counts constitute lower bounds for the NR case. Starting with the set of unique trimers, the set of tetramers is calculated in much the same way. We iterate this procedure up to complexes of size 7. The results are shown in Figure 5A. Truncating the enumeration at this point results in nearly  $10^5$  unique molecular species. Unfortunately, for complexes of size 8 or larger the computational cost of checking for duplicates exceeds current computational resources. Despite this limitation, brute-force enumeration up to size 7 indicates that the cSIN is likely to generate a very large number of possible unique complexes.

*Complexes in Random Acyclic Graphs.* We construct random acyclic interaction graphs (RAGs) with



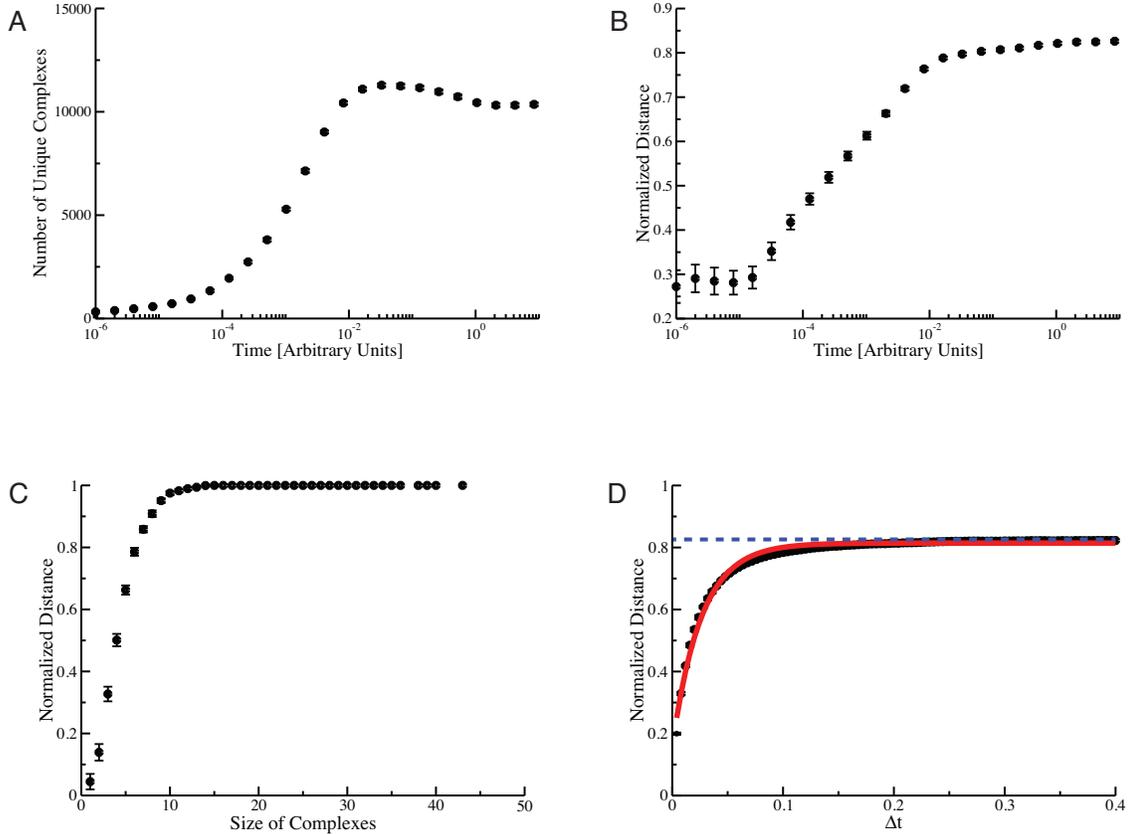
**Figure 5. Combinatorial complexity of the cSIN.** **A:** Panel A reports the number of unique complexes that could be produced by the cSIN as a function of complex size using brute force enumeration. As described in the text, complexes that contain more than one copy of a particular protein are discarded, since they could correspond to polymers. Given that the NR constraint allows for multiple copies of a protein to enter a complex in certain situations (see section 7.1 of the Supporting Information), the numbers displayed here represent a lower bound on the number of unique complexes for the NR constraint. The red line represents an exponential regression of the data, with  $y = 69.6e^{0.89x}$ . **B:** Panel B reports the estimated combinatorial complexity of cSIN-like acyclic networks as a function of network size, using the procedure described in section 3 of the Supporting Information. Each point represents an average over 10 independently generated model networks with the same edge density as the cSIN. The red line depicts an exponential regression with  $y = 2.74e^{0.75x}$

varying number  $N$  of nodes but a fixed cSIN edge density  $\rho \approx 0.039$  and compute the number of possible complexes, as detailed in section 4 of the Supporting Information. Each point in Figure 5B reports the average number from 10 independently generated RAGs with a given  $N$ . Although we cannot give a tight estimate for the cSIN, we conclude from Figure 5B that the number of possible unique cSIN complexes is in the range of  $10^{30}$  to  $10^{40}$ , which is much larger than the total number of proteins present in any given yeast cell. This approach assumes, however, that all possible complexes can be physically realized. In section 5 of the Supporting Information, we describe a simple calculation to estimate the consequences that steric constraints might have on the total number of molecular species that an interaction network could form. The case we considered represents a fairly strong constraint, in which steric effects become more and more prominent as complexes get larger. Given that the surface area of a complex will tend to increase with increasing size, this might not represent the most realistic situation, but the model demonstrates that even strong steric constraints do not curtail combinatorial complexity significantly. If only 20% of complexes of a given size can be realized, the total number is still  $\sim 10^{12}$ , suggesting that steric constraints would have to be incredibly strong in order to reduce the number of molecular possibilities to numbers that allow their simultaneous sampling by a cell.

## Network dynamics with uniform affinities

Based on our assumptions about affinities and rate constants (Methods section), uniform affinities translate into uniform rate parameters. The case we discuss here consists in a stochastic dissociation constant  $\kappa_D = 250$  molecules (corresponding to a deterministic  $K_D = 10$  nM); a stochastic on-rate  $\beta_+ = 0.01$  molecule $^{-1}$  s $^{-1}$  (corresponding to a deterministic on-rate  $2.5 \times 10^8$  M $^{-1}$ s $^{-1}$ ); and a stochastic off-rate  $\beta_- = 2.5$  s $^{-1}$  (corresponding to a deterministic off-rate  $k_- = 2.5$  s $^{-1}$ ). Results for other uniform interac-

tion strengths are similar and are discussed in the Supporting Information.

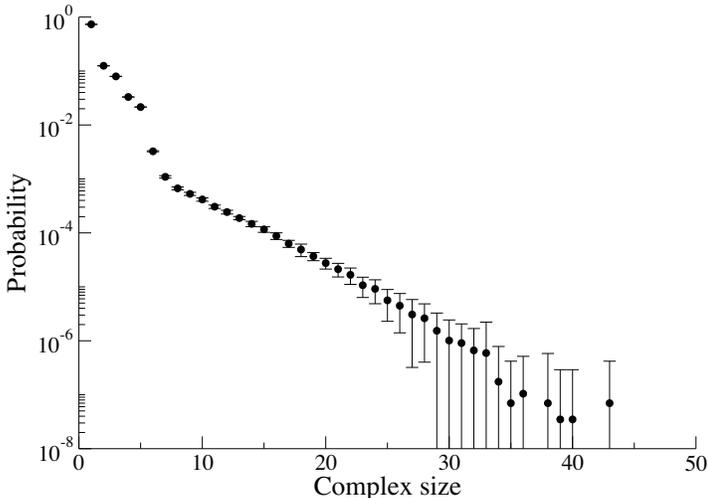


**Figure 6. Dynamic diversity of the cSIN in yeast cells.** **A:** The graph reports the number of unique complexes actually present in a simulated system (“cell”) as a function of time. Each point represents an average over 15 independent simulations. In all panels of this figure, the error bars represent approximately 95% confidence intervals. **B:** The normalized distance between the complement of complexes (“complexomes”) generated by individual simulations is shown as a function of time. Each point is an average over all unique comparisons between 15 independent simulations. Using the parameters described in the text, the separation between steady states reaches  $\sim 80\%$  of the maximal distance. **C:** The stationary distance between cells is shown as a function of complex size, averaged over all of the unique comparisons between 15 independent simulations. The complexomes of cells are nearly identical with regard to small complexes, due to fewer combinatorial possibilities and the high relative abundance of small complexes (see Figure 7 below). However, complexomes differ dramatically for large complexes. This is the case for all combinations of parameters and ring closure scenarios we have tested (see below and the Supporting Information). Since other parameter sets do not substantially change the relationship shown here, much of the difference in inter-cell distances for these parameter sets derives from how heavily the dynamics sample large complexes. **D:** The distance between a cell at time  $t$  and the same cell at time  $t + \Delta t$  is shown as a function of  $\Delta t$ . The first time point  $t$  is taken after cells have reached steady state (in this case,  $t = 2$ , see panels A and B). The blue line denotes the average inter-cell distance at steady state, taken from the last time point in panel A above. The red curve represents an exponential fit to the relaxation, with  $y = 0.81 - 0.66e^{-38x}$ .

The number of unique molecular species present as a function of time (averaged over 15 independent simulations) is shown in Figure 6A. The system approaches a steady-state comprising around 10,000 unique complexes. The approach to steady state occurs on a time scale that corresponds roughly to the equilibration of individual binding reactions. Significantly weaker interactions lead to somewhat fewer unique species, as does the SR scenario. In all cases, no single (simulated) cell contains enough unique complexes to even sample all of the 7-mer structures compatible with the network (Figure 5A), much less the set of all possible complexes. To characterize the differences between simulations, or independent “cells”, we define the set of unique complexes in a cell  $i$  as  $C_i$  and the distance between two cells  $i$  and  $j$  as:

$$d(i, j) = \frac{|C_i \Delta C_j|}{|C_i \cup C_j|} \tag{5}$$

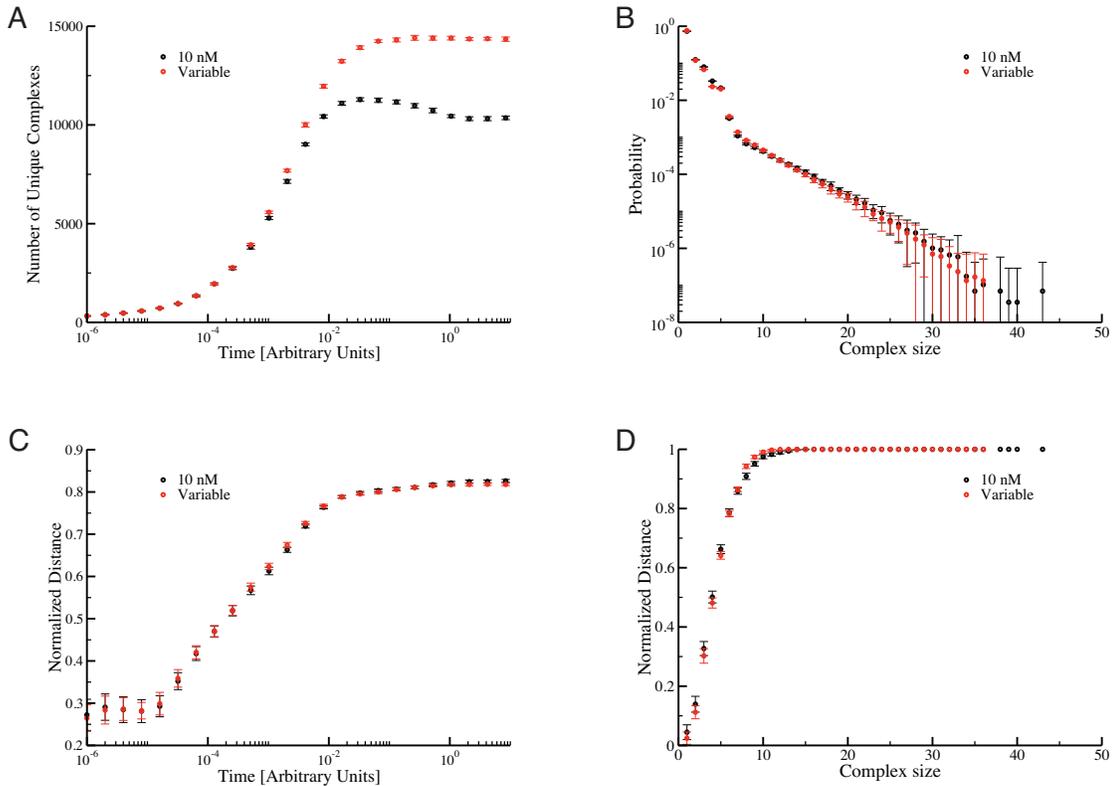
where  $|X|$  denotes the number of elements in set  $X$  and  $C_i \Delta C_j$  denotes the symmetric difference (i.e. the set of complexes that are either in cell  $i$  or cell  $j$ , but not both). Normalizing the symmetric difference by the union  $C_i \cup C_j$  results in a  $d(i, j)$  representing the probability that a particular type of complex found in either cell  $i$  or cell  $j$  is unique to one cell or the other. Although cells start out as identical, they rapidly diverge to a distance of about 0.83, indicating that only 17% of complexes are found in both cells at steady-state (Figure 6B). Alternative distance functions, including definitions that consider differences in copy number, produce similar results (Supporting Information). The exact value of the steady-state distance depends on details and parameters of the simulations: The SR scenario leads to lower distances—as low as  $\sim 0.4$  (Supporting Information).



**Figure 7. Distribution of complex sizes.** The graph shows the distribution of complex sizes for NR simulations with all dissociation constants set to 10 nM. This distribution is calculated at the final time point for the simulations represented in Figure 6. The points on the graph represent the average probability of finding a complex of a certain size across 15 independent simulations. The error bars in this case are set to approximate 95% confidence intervals; for large complexes, the error bars exceed the scale for the lower bound. This is because the 95% confidence intervals include 0, which cannot be displayed on the logarithmic scale of the ordinate.

The divergence of initially identical cells in the space of possible complexes varies strongly with complex size and copy number (Figures 6C and section 8 of the Supporting Information). All cells exhibit an essentially identical repertoire of monomers, dimers and trimers, which tend to be the most

common complexes. However, for complexes of size 9 or larger, cells tend to be completely distinct from one another. We generally find only a single example of any given large complex in a cell, and any particular large complex found at time  $t$  in one cell will not be found anywhere else in the population (Figure 6C). This finding is robust to changes in the affinity parameters and characterizes both the SR and NR constraints (Supporting Information).



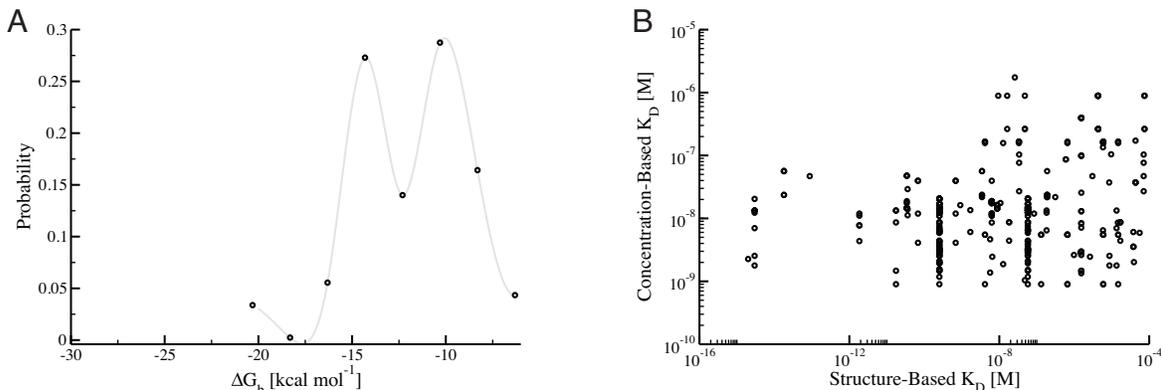
**Figure 8. Comparison between network dynamics based on uniform affinities and concentration-based affinities.** **A:** The number of unique complexes in independent simulations as a function of time: each curve represents the average over 15 independent simulations. In this panel, as with all of the panels in this figure, the error bars represent  $\approx 95\%$  confidence intervals. Allowing interaction strengths to vary across the network produces more unique complexes at steady state ( $\sim 15000$  for the variable case compared to  $\sim 10000$  for the 10 nM case). **B:** Comparison of the distribution of complex sizes: the distributions represent the probability of finding a complex of a particular size across the entire population of 15 simulations at the final time point in panel A. The two interaction affinity scenarios produce similar distributions, with the 10 nM simulations sampling somewhat larger complexes. **C:** Comparison of the distance between independent simulations over time: each curve represents the average over all unique comparisons between 15 independent simulations using the distance measure defined in equation 5. As in panel B, the two scenarios produce essentially identical curves. **D:** Comparison of the distance between independent simulations as a function of complex size: each curve represents the average over all unique comparisons between 15 independent simulations at the final time point in panel A. Again, the two parameter scenarios produce essentially the same result.

Figure 7 shows the distribution of complex sizes at steady state. This distribution is derived from the same set of simulations examined in Figure 6. Small complexes (i.e. monomers and dimers) clearly dominate the distribution, with larger complexes being comparatively rare. The dominance of monomers in this case is somewhat surprising; the interactions here are fairly strong, so one would expect most proteins to participate in at least one complex. The empirical distribution of protein copy numbers, however, is approximately log-normal [40]. The most common protein in these simulations is present with over  $10^5$  copies, while the least common protein has only  $\sim 100$  copies. Thus, certain proteins are present at much higher concentration than any of their potential binding partners, leaving many of the former as monomers. Although quite rare, the largest complexes sampled by these simulations have over 40 members.

These results suggest that each cell on its own might drift in the space of complexes. As seen in Figure 6D, the distance between a particular cell at times  $t$  and  $t + \Delta t$  rapidly increases. For a realistic binding rate ( $\sim 10^7 \text{ s}^{-1}\text{M}^{-1}$ ) [48], the time-scale on which a cell loses memory of its former “compositional self” is  $\sim 0.3$  seconds. We refer to the independent sampling of a distinct and constantly varying set of complexes over time as “compositional drift”.

## Network dynamics with concentration-based affinities

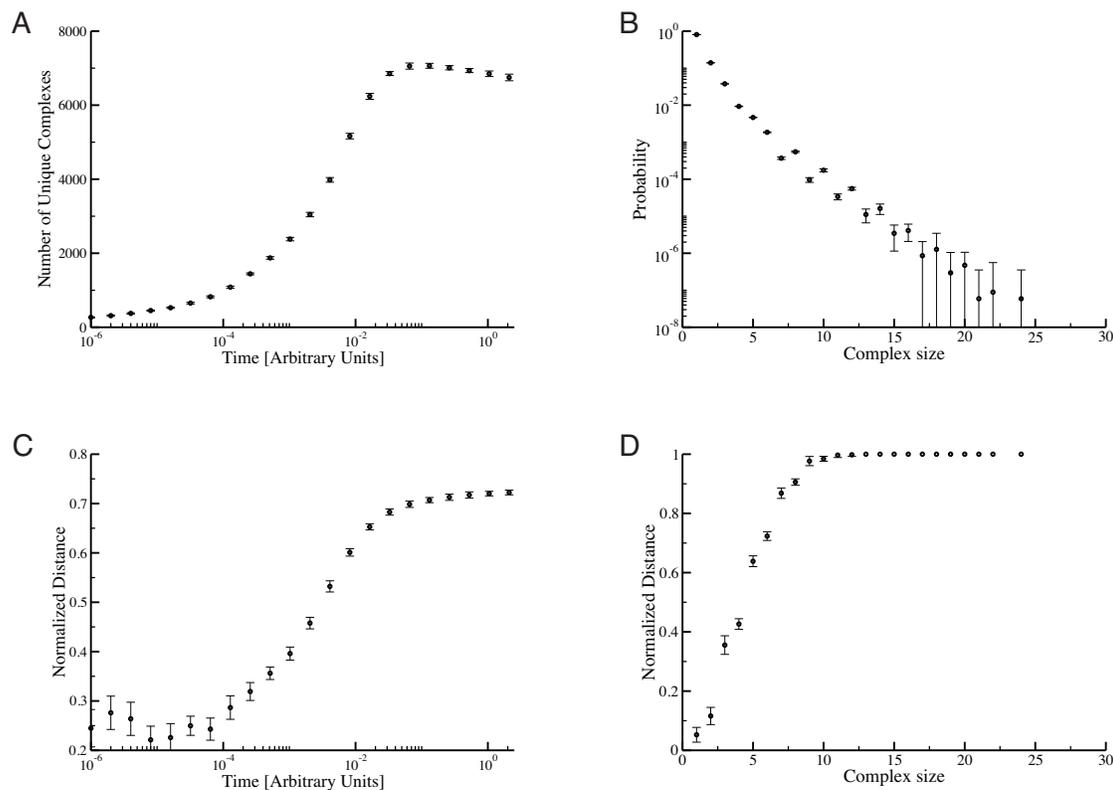
We find that simulations in which  $K_D$ ’s vary across the network according to equation 3 produce results very similar to those obtained at 10 nM for the NR scenario. Figure 8 exhibits the appropriate comparisons. The qualitative results are the same for the SR scenario, with lower affinities leading to somewhat smaller average distances (data not shown) but still large distances for large complexes.



**Figure 9. Binding free energies and dissociation constants for the cSIN2.** **A:** A plot of the distribution of free energies for reactions in the cSIN2. The black circles are a histogram of the free energies; the grey line represents a smoothed version of the distribution. The average free energy is  $-11.0 \text{ kcal mol}^{-1}$ , which corresponds to a dissociation constant of 10.6 nM. **B:** This plot presents a comparison of the structure-based  $K_D$ ’s for each edge in the cSIN2 (abscissa) and the concentration-based  $K_D$ ’s (ordinate). For each interaction in the cSIN2 the concentration-based  $K_D$  is obtained using equation 3. Despite the similarity in the average affinity in both cases (corresponding to a  $K_D$  of around 10 nM), the two methods produce  $K_D$  values that are very different from one another: the linear correlation produces an  $R^2$  of 0.04.

## Network dynamics with structure-based affinities

Proceeding as detailed in “Affinities” of the Methods section, we constructed a version of the cSIN—the cSIN2—in which each binding affinity in the network was calculated from the change in non-polar solvent-accessible surface area based on the protein structures originally used to construct the SIN itself.



**Figure 10. Results from NR simulations of the cSIN2.** **A:** The number of unique complexes in independent simulations as a function of time: this curve represents the average over 15 independent simulations. In this panel, as with all other panels in this figure, the error bars represent  $\approx 95\%$  confidence intervals. The steady-state number of unique complexes is slightly smaller for the cSIN2 than the original cSIN using constant 10 nM affinities ( $\sim 7000$  compared with  $\sim 10000$ ). **B:** This plot shows the probability of finding a complex of a particular size across the entire population of 15 simulations at the final time point in panel A. The distribution of sizes is similar to that found for NR simulations of the original cSIN, although the complexes are, on average, somewhat smaller than those obtained from NR simulations of the cSIN at 10 nM. **C:** This plot displays the distance between independent simulations over time: the curve represents the average over all unique comparisons between 15 independent simulations using the distance measure defined in equation 5. The distances obtained from the cSIN2 are slightly lower than those obtained from the cSIN at 10 nM ( $\sim 0.72$  vs.  $\sim 0.83$ ). **D:** This curve represents the distance between simulations as a function of complex size, averaged over all unique comparisons between 15 independent simulations at the final time point in panel A. The overall shape of this curve is essentially identical to the 10 nM case for the original cSIN as displayed in Figure 5; the main difference is that the simulations based on structure-derived  $K_D$ 's sample somewhat smaller complexes than the original 10 nM case.

The cSIN2 consists of 414 edges between 166 nodes. A number of edges in the original cSIN are lost in constructing the cSIN2, because some domain-domain interactions do not have representative structures in the iPfam database that are truly intermolecular, while others do not have structures where binding is strong enough (see section 8.2 of the Supporting Information). The distribution of free energies of binding,  $\Delta G_b$ , for the cSIN2 is shown in Figure 9A. It has an average of  $-11.0$  kcal mol $^{-1}$  with a standard deviation of 2.96 kcal mol $^{-1}$ . Interestingly, this average free energy corresponds to a dissociation constant of 10.6 nM which is close to the average free energy seen in the PINT database [21] and used for all of the interactions in the simulations described above under the uniform rate constant scenario.

The concentration-based  $K_D$  scenario (i.e. the case in which dissociation constants are derived from equation 3) yields an average affinity that is very similar to the structure-based  $K_D$ 's ( $K_D$ 's of 13.1 and 10.6 nM, respectively). However, despite the similarity in the average, the  $K_D$  values for the structure-based affinities vary considerably across the network in a manner that appears independent from the concentration-based affinities derived from equation 3, Figure 9B.

Figure 10 summarizes the results of NR simulations of the cSIN2 using these structure-based affinities. As can be seen from Figure 10, the overall behavior of the cSIN2 is very similar to that of the original cSIN simulated with NR constraints. The cSIN2 yields somewhat lower steady-state distances than the original cSIN when simulated using 10 nM affinities ( $\sim 0.72$  vs.  $\sim 0.83$ ) or 100 nM affinities (Supporting Information), largely because the cSIN2 simulations sample somewhat fewer large complexes. SR simulations based on the cSIN2 are also very similar to the 10 nM SR case (data not shown).

## Other results

The Supporting Information includes discussions of simulations using alternative distance measures (equation 5); comparisons between different uniform affinities; and the global SR scenario. The thermodynamics of ring-like protein complexes (discussed in section 6.1 of the Supporting Information) can give rise to situations in which a particular pair of sites might not bind one another strongly enough to be detected in a high-throughput interaction screen (such as a Yeast Two-Hybrid experiment) but could nonetheless contribute dramatically to the stability of certain complexes by forming a bond to complete a ring. In the Supporting Information we discuss the addition of such "cryptic cycles". All these variations leave the main observation of compositional drift intact.

## Discussion

Our simulations provide a dynamical picture of PPI networks based on a model that is respectful of their combinatorial complexity. PPI networks represent binding capabilities between proteins typically determined by an assay that yields inherently local information. Two broad components were necessary for making and running a model of a PPI network: (i) A representation of the system that can handle combinatorial complexity implicitly, since the number of possible complexes is astronomical, preventing their explicit representation. (ii) A dataset in which the interactions derived from a binding assay have been curated, and binding interactions are resolved at the level of domains or sites, allowing the distinction between interactions that are mutually compatible and those that are mutually exclusive. The first component is addressed by rule-based approaches, such as Kappa or BNGL. The second component is a suitable dataset that has been recently compiled by Kim et al [11]. We bring these two critical components together, along with protein localization, abundance data and a few biophysical assumptions, to generate a simulation of a large slice of a PPI network.

According to our simulations, systems that start from identical initial conditions diverge from one another rapidly with regard to the complexes they contain, eventually sampling different regions of the space of possible complexes. This is particularly the case for large complexes, where independent simulations tend to be essentially disjoint. Our model indicates that the complexity of such networks will

result in compositional drift, even with the biophysical constraints imposed by the NR and SR scenarios. However, we consider neither post-translational modifications nor translation and degradation processes. Our systems therefore reach thermodynamic equilibrium. At equilibrium the vast space of molecular possibilities permits energetically neutral compositional drift, i.e. a never-ending change in the set of realized complexes present in a particular simulation.

The data from which our network is built has clear limitations. High-throughput methods for acquiring PPI data, such as Y2H assays, tend to have substantial false positive and false negative rates [11, 42, 50]. Curated, structure-based data sets like the SIN alleviate this drawback to some extent, but we cannot rule out the presence of fictitious edges in the cSIN network. Given that drift, especially among large complexes, is a robust feature of our simulations, it is unlikely that the ultimate removal of such edges would affect this phenomenon. Indeed, the cSIN2, which contains a slightly smaller set of interactions based on more stringent structural evidence, undergoes essentially the same level of drift as other versions of the network, indicating that inaccuracies in the underlying interaction data are unlikely to have a large influence on the overall dynamics described here (although they would have an influence on the identity of the complexes formed).

Our dynamic model does not include synthesis and degradation processes, raising the question whether limiting the time proteins persist in the cell might affect drift. High-throughput measurements of protein degradation rates [51] indicate that the average half-life of yeast proteins is around 42 minutes, with a minimum observed half-life of about 2 minutes. In our simulations, both the total number of unique complexes and their size distribution generally reach equilibrium in about one second (see, e.g., Figure 6A). Degradation processes are thus unlikely to occur at high enough rates to fundamentally influence the average size of complexes at steady-state and thus the presence of drift. However, in the SR scenario, ring-like structures are by definition so stable that they are much more likely to be removed by degradation or dilution than spontaneous dissociation. In that case, it is conceivable that degradation actually increases drift on longer timescales. Given our current computational limitations, we are unable to carry out simulations that are long enough to assess the influence of realistic synthesis and degradation rates on drift in the SR scenario.

The empirical data that define our model are also too limited and fragmentary to provide an accurate reflection of the actual geometric, kinetic, and biological constraints that determine complex formation. Indeed, large molecular machines like the ribosome and the proteasome are highly unlikely to undergo compositional drift [52–54]. In view of these shortcomings, what are we to make of compositional drift? At a conceptual level, our work suggests a serious problem that must be overcome in order for such complexes to assemble reliably in the cell. It is not enough for the parts of a specific supra-molecular complex to simply “fit together snugly” or bind with high affinity when independent binding sites and a large number of extraneous binding partners yield a fantastically large set of combinational possibilities that can never be exhaustively populated. Absent any further constraints, the system becomes “lost” in the vast set of possible species available to it, preventing the reliable assembly of a desired target complex.

The reduction of drift requires limiting the space of possibilities available to a PPI system. One strategy to accomplish this would be to limit the size of complexes that can form, since small complexes are well-sampled in our simulations and do not exhibit significant drift. A second strategy would be to evolve “hierarchical” assembly pathways, thus curtailing the number of accessible complexes but not necessarily their size. A simple implementation of the first strategy would be to constrain the number of sites in proteins, especially those proteins that are “hubs” in the network. Such an architecture resembles the scenario depicted at the top of Figure 1, but it does not seem to characterize the overall SIN or the cSIN studied here. Moreover, such a network architecture would not account for large macromolecular machines. A flexible implementation of the second strategy is the use of conditional rules, where binding interactions between sites are highly sensitive to the molecular context in which they occur. There are many potential mechanisms suitable for introducing causal dependencies between binding and unbinding events: for instance, allostery and cooperativity could be employed to radically alter the binding free en-

ergy of a particular interaction in specific contexts, thus inducing the dynamics to avoid a large fraction of molecular possibilities. Post-translational modifications could also be used to create causal dependencies, provided they are deployed in such a manner as not to increase the combinatorial complexity [55].

We view compositional drift as the network analogue of the protein folding *problem*. The combinatorial explosion of possible conformational states available to the polypeptide chain raised the conundrum of how a protein can fold quickly and stably into a native structure (the so-called “Levinthal paradox”). The exploration of this problem eventually led to a framework for identifying the evolved features of free energy landscapes that ensure reliable folding of proteins [56, 57]. Likewise, the combinatorial explosion of possible molecular associations gives rise to the compositional drift problem for assembly in a network context. While there are many potential mechanisms suitable for introducing causal dependencies between binding and unbinding events, the specific deployment of these mechanisms can only be understood in light of the system-wide drift problem that they solve. In other words, compositional drift brings to light the need for complex networks to evolve particular *chemical potential landscapes* in order for assembly to proceed reliably within cells. This also raises the question, especially with regard to the many transient protein associations that can be formed during signaling, whether it is at all possible to entirely eliminate drift while reusing proteins in diverse contexts within the same cell. A certain level of compositional drift might be unavoidable, and in some situations could actually constitute an evolutionarily advantageous source of non-genetic individuality in isogenic populations.

## Acknowledgments

The authors would like to thank Drs. Javier Apfeld, Russ Harmer, Tom Kolokotronis, Sergei Maslov, and Ethan Perlstein for their comments on the manuscript. E.J.D. was partially supported by an NRSA postdoctoral fellowship from the NIH.

## References

1. Uetz P, Giot L, Cagney G, Mansfield TA, Judson RS, et al. (2000) A comprehensive analysis of protein-protein interactions in *saccharomyces cerevisiae*. *Nature* 403: 623–627.
2. Ito T, Chiba T, Ozawa R, Yoshida M, Hattori M, et al. (2001) A comprehensive two-hybrid analysis to explore the yeast protein interactome. *Proceedings of the National Academy of Sciences of the United States of America* 98: 4569–4574.
3. Gavin AC, Aloy P, Grandi P, Krause R, Boesche M, et al. (2006) Proteome survey reveals modularity of the yeast cell machinery. *Nature* 440: 631–636.
4. Yu H, Braun P, Yildirim MA, Lemmens I, Venkatesan K, et al. (2008) High-quality binary protein interaction map of the yeast interactome network. *Science* 322: 104–110.
5. Jones RB, Gordus A, Krall JA, MacBeath G (2006) A quantitative protein interaction network for the *erbB* receptors using protein microarrays. *Nature* 439: 168–174.
6. Stellberger T, Häuser R, Baiker A, Pothineni VR, Haas J, et al. (2010) Improving the yeast two-hybrid system with permutated fusion proteins: the Varicella Zoster Virus interactome. *Proteome science* 8: 8.
7. Oda K, Matsuoka Y, Funahashi A, Kitano H (2005) A comprehensive pathway map of epidermal growth factor receptor signaling. *Molecular Systems Biology* 1: 2005.0010.
8. Csete M, Doyle J (2004) Bow ties, metabolism and disease. *Trends in Biotechnology* 22: 446–450.

9. Jeong H, Tombor B, Albert R, Oltvai ZN, Barabasi AL (2000) The large-scale organization of metabolic networks. *Nature* 407: 651–654.
10. Goldberg DS, Roth FP (2003) Assessing experimentally derived interactions in a small world. *Proc Natl Acad Sci U S A* 100: 4372–4376.
11. Kim PM, Lu LJ, Xia Y, Gerstein MB (2006) Relating three-dimensional structures to protein networks provides evolutionary insights. *Science* 314: 1938–41.
12. Hlavacek W, Faeder J, Blinov M, Perelson A, Goldstein B (2003) The complexity of complexes in signal transduction. *Biotechnol Bioeng* 84: 783–794.
13. Hlavacek WS, Faeder JR, Blinov ML, Posner RG, Hucka M, et al. (2006) Rules for modeling signal-transduction systems. *Science STKE* 344: re6.
14. Danos V, Laneve C (2004) Formal molecular biology. *Theoretical Computer Science* 325: 69–110.
15. Blinov ML, Faeder JR, Hlavacek WS (2004) BioNetGen: Software for rule-based modeling of signal transduction based on the interactions of molecular domains. *Bioinformatics* 20: 3289–3292.
16. Blinov ML, Yang J, Faeder JR, Hlavacek WS (2006) Graph theory for rule-based modeling of biochemical networks. *Lect Notes Comput Sci* 4230: 89–106.
17. Danos V, Feret J, Fontana W, Harmer R, Krivine J (2007) Rule-based modelling of cellular signalling. In: *Proceedings of the 18th Int. Conf. on Concurrency Theory*. Lisboa, Portugal: Springer, volume 4703 of *Lecture Notes in Computer Science*, pp. 17–41.
18. Danos V, Feret J, Fontana W, Krivine J (2007) Scalable simulation of cellular signalling networks. In: *Proceedings APLAS 2007*. Springer, volume 4807 of *Lecture Notes in Computer Science*, pp. 139–157.
19. Yang J, Monine MI, Faeder JR, Hlavacek WS (2008) Kinetic monte carlo method for rule-based modeling of biochemical networks. *Phys Rev E* 78: 031910.
20. Faeder JR, Blinov ML, Hlavacek WS (2009) Rule-based modeling of biochemical systems with bionetgen. *Methods Mol Biol* 500: 113–67.
21. Maslov S, Ispolatov I (2007) Propagation of large concentration changes in reversible protein-binding networks. *Proc Natl Acad Sci U S A* 104: 13655–13660.
22. Jeong H, Mason S, Barabasi A, Oltvai Z (2001) Lethality and centrality in protein networks. *Nature* 411: 41–42.
23. Thomas A, Cannings R, Monk N, Cannings C (2003) On the structure of protein-protein interaction networks. *Biochemical Society transactions* 31: 1491–1496.
24. Barabási AL, Oltvai ZN (2004) Network biology: understanding the cell’s functional organization. *Nature Reviews Genetics* 5: 101–113.
25. Zotenko E, Mestre J, O’Leary DP, Przytycka TM (2008) Why do hubs in the yeast protein interaction network tend to be essential: reexamining the connection between the network topology and essentiality. *PLoS Computational Biology* 4: e1000140.
26. Ghaemmaghami S, Huh WK, Bower K, Howson RW, Belle A, et al. (2003) Global analysis of protein expression in yeast. *Nature* 425: 737–41.

27. Huh WK, Falvo JV, Gerke LC, Carroll AS, Howson RW, et al. (2003) Global analysis of protein localization in budding yeast. *Nature* 425: 686–91.
28. Perktold A, Zechmann B, Daum G, Zellnig G (2007) Organelle association visualized by three-dimensional ultrastructural imaging of the yeast cell. *FEMS Yeast Res* 7: 629–38.
29. Cho RJ, Campbell MJ, Winzeler EA, Steinmetz L, Conway A, et al. (1998) Modeling networks of coupled enzymatic reactions using the total quasi-steady state approximation. *Mol Cell* 2: 65-73.
30. Spellman PT, Sherlock G, Zhang MQ, Iyer VR, Anders K, et al. (1998) Comprehensive identification of cell cycle-regulated genes of the yeast *Saccharomyces cerevisiae* by microarray hybridization. *Mol Biol Cell* 9: 3273–3297.
31. Shedden K, Cooper S (2002) Analysis of cell-cycle gene expression in *Saccharomyces cerevisiae* using microarrays and multiple synchronization methods. *Nucl Acids Res* 30: 2920–2929.
32. de Lichtenberg U, Jensen LJ, Fausboll A, Jensen TS, Bork P, et al. (2005) Comparison of computational methods for the identification of cell cycle-regulated genes. *Bioinformatics* 21: 1164–1171.
33. de Lichtenberg U, Jensen LJ, Brunak S, Bork P (2005) Dynamic Complex Formation During the Yeast Cell Cycle. *Science* 307: 724–727.
34. Danos V, Feret J, Fontana W, Harmer R, Krivine J (2008) Rule-based modelling, symmetries, refinements. In: *Formal Methods in Systems Biology*. Cambridge, UK: Springer, volume 5054 of *Lecture Notes in Bioinformatics*, pp. 103-122.
35. Danos V, Feret J, Fontana W, Krivine J (2008) Abstract interpretation of cellular signalling networks. In: *Verification, Model Checking, and Abstract Interpretation*. Springer, volume 4905 of *Lecture Notes in Computer Science*, pp. 83–97.
36. Danos V, Feret J, Fontana W, Harmer R, Krivine J (2009) Rule-based modelling and model perturbation. *Transactions on Computational Systems Biology* 11: 116-137.
37. Harmer R, Danos V, Feret J, Krivine J, Fontana W (2010) Intrinsic information carriers in combinatorial dynamical systems. *Chaos* 20: 037108.
38. Gillespie DT (1976) A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of Computational Physics* 22: 403–434.
39. Kumar MD, Gromiha MM (2006) Pint: Protein-protein interactions thermodynamic database. *Nucleic Acids Res* 34: D195-8.
40. Ghaemmaghami S, Huh WK, Bower K, Howson RW, Belle A, et al. (2003) Global analysis of protein expression in yeast. *Nature* 425: 737-41.
41. Zhang J, Maslov S, Shakhnovich EI (2008) Constraints imposed by non-functional protein-protein interactions on gene expression and proteome size. *Mol Syst Biol* 4: 210.
42. Deeds EJ, Ashenberg O, Shakhnovich EI (2006) A simple physical model for scaling in protein-protein interaction networks. *Proc Natl Acad Sci U S A* 103: 311–316.
43. Horton N, Lewis M (1992) Calculation of the free energy of association for protein complexes. *Protein Sci* 1: 169-81.

44. Bougouffa S, Warwicker J (2008) Volume-based solvation models out-perform area-based models in combined studies of wild-type and mutated protein-protein interfaces. *BMC Bioinformatics* 9: 448.
45. Fraternali F, Cavallo L (2002) Parameter optimized surfaces (pops): analysis of key interactions and conformational changes in the ribosome. *Nucleic Acids Res* 30: 2950-2960.
46. Pang X, Qin S, Zhou HX (2011) Rationalizing 5000-fold differences in receptor-binding rate constants of four cytokines. *Biophys J* 101: 1175–1183.
47. Jorgensen P, Nishikawa JL, Breikreutz BJ, Tyers M (2002) Systematic identification of pathways that couple cell growth and division in yeast. *Science* 297: 395–400.
48. Camacho CJ, Kimura SR, DeLisi C, Vajda S (2000) Kinetics of desolvation-mediated protein-protein binding. *Biophys J* 78: 1094–1105.
49. Saiz L, Vilar JM (2006) Stochastic dynamics of macromolecular-assembly networks. *Mol Syst Biol* 2: 2006 0024.
50. Kuchaiev O, Raajski M, Higham DJ, Prulj N (2009) Geometric de-noising of protein-protein interaction networks. *PLoS Comput Biol* 5: e1000454.
51. Belle A, Tanay A, Bitincka L, Shamir R, O’Shea EK (2006) Quantification of protein half-lives in the budding yeast proteome. *Proc Natl Acad Sci USA* 103: 13004–9.
52. Ban N, Nissen P, Hansen J, Moore PB, Steitz TA (2000) The complete atomic structure of the large ribosomal subunit at 2.4 Å resolution. *Science* 289: 905–20.
53. Yusupov MM, Yusupova GZ, Baucom A, Lieberman K, Earnest TN, et al. (2001) Crystal structure of the ribosome at 5.5 Å resolution. *Science* 292: 883–96.
54. Murata S, Yashiroda H, Tanaka K (2009) Molecular mechanisms of proteasome assembly. *Nat Rev Mol Cell Biol* 10: 104–115.
55. Mayer BJ, Blinov ML, Loew LM (2009) Molecular machines or pleiomorphic ensembles: signaling complexes revisited. *Journal of Biology* 8: 81.
56. Shakhnovich E (2006) Protein folding thermodynamics and dynamics: where physics, chemistry, and biology meet. *Chem Rev* 106: 1559–1588.
57. Onuchic JN, Wolynes PG (2004) Theory of protein folding. *Curr Opin Struct Biol* 14: 70–75.

# Combinatorial complexity and compositional drift in protein interaction networks

## Supporting Information

Eric J. Deeds<sup>1</sup>, Jean Krivine<sup>2</sup>, Jérôme Feret<sup>3</sup>, Vincent Danos<sup>4</sup> and Walter Fontana<sup>5</sup>

<sup>1</sup>Center for Bioinformatics and Department of Molecular Biosciences, The University of Kansas, Lawrence KS 66047, USA

<sup>2</sup>Laboratoire PPS de l'Université Paris 7 and CNRS, F-75230 Paris Cedex 13, France

<sup>3</sup>Laboratoire d'Informatique de l'École normale supérieure, INRIA, ÉNS, and CNRS, 45 rue d'Ulm, F-75230 Paris Cedex 05, France

<sup>4</sup>School of Informatics, University of Edinburgh, Edinburgh, UK

<sup>5</sup>Department of Systems Biology, Harvard Medical School, 200 Longwood Avenue, Boston MA 02115, USA

Email: Eric Deeds - deeds@ku.edu; Jean Krivine - jkrivine@pps.jussieu.fr; Jerome Feret - feret@ens.fr; Vincent Danos - vdanos@inf.ed.ac.uk; Walter Fontana - walter@hms.harvard.edu;

## Contents

<b>1</b>	<b>A Rule-Based Modeling Framework</b>	<b>2</b>
1.1	Kappa . . . . .	2
1.2	Site graphs, contact maps, complexes, and molecular species . . . . .	5
1.3	Locality of rules and cyclical structures . . . . .	6
<b>2</b>	<b>Simulating Kappa-models</b>	<b>7</b>
2.1	An overview of the stochastic simulation method . . . . .	7
2.2	Time advance with null events . . . . .	8
<b>3</b>	<b>Counting complexes in acyclic contact maps</b>	<b>10</b>
<b>4</b>	<b>Random Acyclic Graphs</b>	<b>12</b>
<b>5</b>	<b>The effect of size constraints on the number of possible complexes</b>	<b>13</b>
<b>6</b>	<b>The “stable rings” scenario</b>	<b>14</b>
6.1	The thermodynamic rationale for the “stable rings” scenario . . . . .	14

6.2	The implementation of the “stable rings” scenario . . . . .	16
<b>7</b>	<b>The “no rings” scenario</b>	<b>18</b>
7.1	The implementation of the “no rings” scenario . . . . .	18
7.2	The relationship between “stable rings” and the “no rings” scenario . . . . .	20
<b>8</b>	<b>Additional Results</b>	<b>20</b>
8.1	Alternative definitions of distance . . . . .	20
8.2	Structure-based affinities . . . . .	25
8.3	Results for SR simulations . . . . .	27
8.4	Results for different affinity scenarios . . . . .	30
8.5	Results based on adding “cryptic” cycles . . . . .	30
<b>9</b>	<b>Comparison with Affinity Purification / Mass Spectrometry data</b>	<b>32</b>
	<b>References</b>	<b>34</b>

## 1 A Rule-Based Modeling Framework

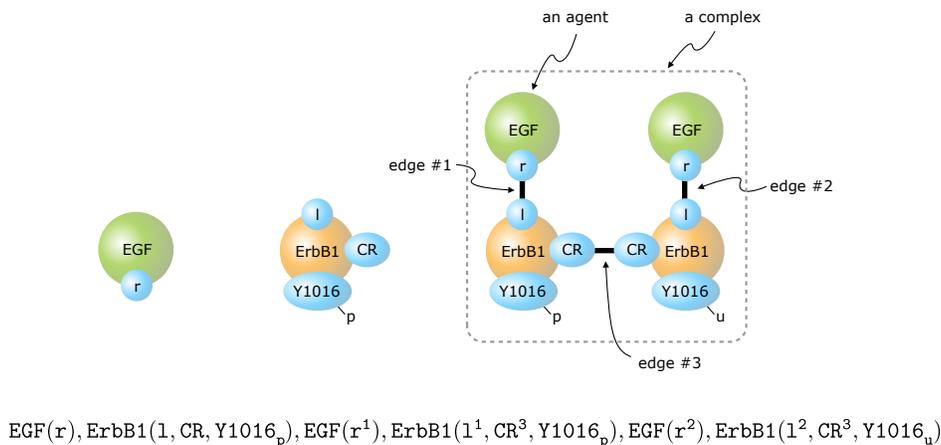
### 1.1 Kappa

Rule-based representations [1–5] are promising for modeling systems based on compositional interactions that can give rise to astronomic numbers of possible molecular species [6]. To model the cytoplasmic SIN (cSIN), we employ a formal language called “Kappa” [2, 4, 7]. Although a formal definition of the language may be found elsewhere (for example, [7, 8]), we include a brief informal description for the sake of a self-contained presentation.

Kappa is a language designed to express interactions between protein agents in terms of rules that refer to partially decontextualized domains or “sites”, much like rules of organic chemical reactions refer to partially decontextualized functional groups. Most of what we describe next also holds for other rule-based approaches, such as BNGL [1].

*Agents and complexes.* Agents are the atoms of the language and represent not full protein complexes but their individual protein constituents. Specification of an agent requires a name and a set of labeled sites—the interface of the agent. A site can have internal states that might represent post-translational modifications. In addition, a site may be bound to at most one site of another agent to form a complex, as sketched in Figure 1. Sites are best understood as representing resources for interactions (which are specified by rules, as explained below). Kappa aims at representing actions and their dependencies on state, without directly representing the structural underpinnings that make such actions physically possible. In other words, Kappa is meant to represent high-level mechanistic knowledge in a manner that enables the study of *process*, i.e. suites of events whose occurrence enables further events.

A complex is a Kappa expression in which the states of the constituent protein agents need not be fully specified, while a *molecular species* is a complex in which each agent occurs with its full complement of sites in definite states. In other words, a Kappa expression is a pattern representing the set of species that match it; the former is an *intensional* object, the latter is an *extensional* one. The combinatorial explosion is tamed by replacing extensional lists by intensional expressions.



**Figure 1. Kappa expressions.** Bottom: The textual representation of a small reaction mixture containing 6 agents that make up three complexes, identified graphically above. The two complexes on the left are simple agents, while the complex on the right consists of 4 agents connected by pairs of identical superscripts at the corresponding binding sites. Top: An equivalent graphical representation. Agents and sites (blue) are labelled with names. An internal state, such as the phosphorylated state  $p$  of site Y1016 at agent ErbB1, is indicated by a labeled barb attached to a site. A mixture of molecular agents is a “site graph”, i.e. a graph whose nodes are sets of sites (each set representing a particular agent) and whose edges are anchored by sites. A complex is a connected site graph, in which a site can bear at most one edge. The graphs do not represent geometric properties; they only convey connectivity and state information.

*The meaning of “agent-based”.* The term “agent-based” has multiple meanings. One meaning refers to any representation of a system based on discrete entities (particles), typically for the purpose of stochastic simulation. Our use of the concept is more nuanced, as it refers to a level of structural granularity, not just discreteness. An analogy to chemistry might help. If molecules are considered to be agents, then the representation of molecular systems requires declaring an infinite alphabet—one symbol (a proper name) per molecular species. In contrast, if atoms are agents, then a very small alphabet and a few rules of *grammar* suffice to build (and thus compositionally name) an infinity of molecules. In both cases we deal with discrete entities and their interactions, but only the latter hinges on a structured language.

*Rules.* Sites represent capabilities for interaction, like binding and post-translational modification, specified by rules. The idea of a rule is to stipulate only the context required for an interaction (Figure 3 of the main text), along with rate information. Instead of directly writing reactions between exhaustively specified molecular species, we write rules that mention names of

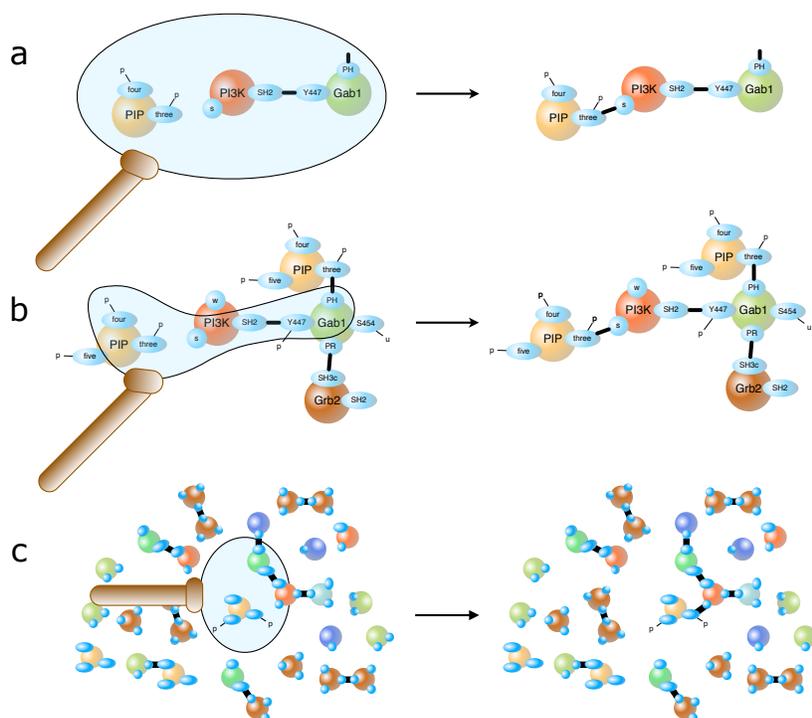
protein-agents and some, but not necessarily all, of their respective sites. In this way, a rule need only make explicit those aspects that are relevant to the interaction being described. More specifically, the left hand side (LHS) of a rule is a pattern. The right hand side (RHS) defines the changes that occur when the LHS is matched in a mixture of agents. The difference between RHS and LHS is called the action of the rule. Sites mentioned on the LHS are said to be tested by the rule. Sites that are tested but not modified constitute the context of a rule’s action. Because rules typically do not mention all the sites and states of an agent, they keep combinatorial complexity implicit. Yet, all possibilities are still realizable at the level of agents in the mixture to which the rules apply (see below). For example, if a substrate has four independent phosphorylatable sites, only one rule is needed to express the modification of any given site, whereas 8 reactions are necessary to express that same action in all possible contexts occasioned by the remaining three sites. Rules can be modified, refined or coarsened as needed to express new knowledge and new hypotheses.

*Events.* Rules are applied to a *mixture*, which is a large Kappa expression representing the contents of a reaction system at a particular time, as illustrated in Figure 2. The lens over the LHS pattern is meant to suggest that applying a rule requires searching for a match, which then creates a reaction instance between the matching complexes (Figure 2B). To apply a rule then means updating the states of the agents so identified in the mixture. Rules may modify several agents at once by invoking any of five elementary actions defined in Kappa: binding, unbinding, internal state modification, creation of an agent, and removal of an agent. In our specific study of complex formation, only binding and unbinding actions occur.

*Models.* A Kappa-model is a collection of rules and an initial mixture on which the rules act. At any given time, different rules may apply in distinct ways involving distinct reactant combinations. The dynamics generated by a Kappa model is therefore stochastic. A rule fires with a probability computed from the number of its matchings in the mixture (mass action) and its rate constant (section 2.1). The overall procedure for generating trajectories that are probabilistically compatible with the underlying master equation of stochastic chemical kinetics follows the Doob-Gillespie algorithm [9, 10] appropriately generalized to rules [11, 12].

*Textual notation.* The graphical rendering of Kappa expressions is equivalent to a textual notation, Figure 1. Although the graphical depiction appeals to instant visual comprehension, it is less convenient for in-text use.

*Implementation and availability.* The Kappa framework is an evolving suite of tools designed for navigating and exploring the structure and dynamics of complex signaling systems. It includes a graphical user interface, a scalable stochastic simulator, a set of procedures for analyzing the causal relationships between rules, a sampler of causal flows (pathways), and an exact coarse-graining procedure for converting a rule-based model into a corresponding set of differential equations whose variables refer to *sets* of molecular species that are distinguishable by the system. The framework and its source codes are available at [www.kappalanguage.org](http://www.kappalanguage.org).



**Figure 2. Rules, instances, and events in Kappa.** **A:** A (fictitious) Kappa-rule specifying conditions for a binding event between two proteins PI3K and PIP. The lens symbolizes the pattern that must be matched for the rule to fire. **B:** A particular combination of complexes is seen to match the LHS of the rule in panel A, yielding a reaction instance. **C:** An event is a particular occurrence of a reaction instance, involving specific molecules in a mixture. Many matchings may be possible for any given rule in a mixture, and many different rules may be applicable at a given moment.

## 1.2 Site graphs, contact maps, complexes, and molecular species

*Site graph.* A site graph is a graph in which nodes are *sets* of sites with edges connecting sites, Figure 1 of the main text. We can think of a site graph as adding additional structure on top of a standard graph by partitioning its nodes into sets. These sets then become the nodes of the site graph and the former nodes—now called sites—remain the endpoints of edges.

*Contact map.* A contact map is a site graph that summarizes statically the possible interactions given by the rules of a Kappa model. The nodes of the contact map are the agents (proteins) that occur in the model. An edge is placed between sites of two agents if the model contains a binding rule whose LHS can be satisfied given an initial condition and the other rules of the model. (The contact map has therefore a semantic content that requires, in general, a reachability computation [13].) Since a contact map summarizes possibilities, each agent type occurs exactly once, but a site may have more than one incident edge. The cSIN in Figure 2 of the main text is a contact map.

*Complexes and molecular species.* A complex of agents is a site graph (e.g. Figure 1) that is realizable given a contact map. A complex is intended to be a partial description of a molecular species that is compatible with the contact map of the model. As such, a complex may contain more than one occurrence of a given agent type, but every site can have at most one bond. A molecular species is a complex that mentions the full complement of sites (and their states) for each agent.

*Cycles.* Cycles in the contact map have different meaning depending on whether the cyclical path of bonds touches the same site twice. If it does, the cycle in the contact map cannot give rise to a cyclical realization of it (i.e. a complex that is a ring), because in a complex a site can be used at most once. A “proper cycle” in the contact map is a cyclical path of bonds, starting and ending at the same node, which does not touch the same site twice. Such a cycle can give rise to infinite realizations of different complexes: each time the addition of a bond comes full circle in the contact map, the bond need not be to an agent already in the complex, but could be to a new copy, thereby creating polymers of any length.

### 1.3 Locality of rules and cyclical structures

Kappa rules are *local* in the sense that their context is self-contained: verifying that a combination of molecules in the mixture satisfies the LHS of the rule never requires information outside of the LHS. This has subtle consequences.

Consider the rule in Figure 3A of the main text. This rule, call it R, stipulates a binding action between an agent of type A with a free site  $s$  (in any internal state) and an agent of type B with a free site  $p$  (in any internal state). Agents A and B may have other sites, but their states do not figure in the context relevant to R. The fact that A and B are disconnected on the rule’s LHS does not imply that their embeddings (matchings) into the mixture must end up disconnected, as shown in Figure 3B of the main text. Indeed, a condition that requires A and B to belong to different molecules is *inherently non-local*: Given two matching agents of type A and B in the mixture, a verification of the condition would require a scan of all (bond) paths originating at the matching agents to confirm that none of them ends up connecting A and B into a single molecule. This task has a worst-case computational cost bounded by the size of the mixture, not the LHS pattern of the rule. In the same vein, expressing in a local fashion that two agents are connected, requires specifying at least one connecting path. Without providing such a path, the condition that “A and B belong to the same molecule” is non-local, since identifying a connecting path requires a worst-case scenario of scanning of the whole mixture. Thus, “local” means that the size of the effective context of any action does not scale with the size of the system.

On the practical side, locality enables the development of efficient scalable simulators [11] and of static analysis tools [4, 7, 13] that greatly facilitate the process of modeling while also providing insight into possible behaviors of a rule collection. On the conceptual side, locality has a strong physical appeal, as the intuition behind the notion of “mechanism” is that of an action that depends on locally accessible information. While Kappa captures a meaningful level of process and attendant abstract locality, it is too simple to capture aspects of *physical* locality that must constrain our simulations of the cSIN, even though the cSIN data contain no information about the particular geometry of a putative complex.

As illustrated in Figure 3 of the main text, one consequence of locality in Kappa is a potential mismatch between the arity of a rule (the number of connected components on the LHS of a rule) and the actual molecularity of the reactions that the rule induces in the mixture. The mismatch is potential, because whether it actually occurs depends on the initial mixture and the molecular species derivable from it by repeated rule applications. Binary Kappa rules are therefore annotated with two stochastic rate constants: a first-order, volume-independent rate constant and a second-order rate constant with a reciprocal dependence on the reaction volume. In the case of the cSIN, we automatically avoid arity mismatches as a by-product of addressing another consequence of locality: polymerization.

Whenever a binary rule can induce both uni- and bimolecular reactions (Figure 3B of the main text), the unimolecular scenario leads to the closure of cycles, whereas the bimolecular scenario can lead to the formation of polymers. In Figure 3B, the second B-agent that has been picked up in reaction 1 can subsequently bind to another C, which can bind a further A, and so on. Whenever cyclical structures can form in a local rule set, polymers can form too. Without proper constraints, mass action would end up favoring polymerization over ring closure. As can be seen from the contact map of the cSIN, Figure 2 of the main text, the potential for cyclical complexes is large and polymerization would be therefore rampant.

## 2 Simulating Kappa-models

### 2.1 An overview of the stochastic simulation method

The stochastic simulator developed for Kappa follows the logic of the well-known Doob-Gillespie procedure [9, 10], but is implemented in a manner that makes the core loop of the procedure scalable, i.e. independent of the size of the system and the number of possible molecular species that are implied by the rules, as detailed in [11].

In the following, recall that molecular species (complexes) are not represented as opaque units, but rather in terms of their constituent protein agents. Thus, as far as the cSIN is concerned, our virtual cytoplasm consists of about three million protein agents (see main text), explicitly represented in computer memory one-by-one. The state of the cytoplasm at any particular time is a graph over this large set of protein agents, with individual complexes being subgraphs whose link structure is continually modified as the simulation proceeds in accordance with the rules of the cSIN and the NR and SR constraints. (The NR and SR constraints are detailed in sections 7 and 6, respectively, of this document.) At no point is there a counter keeping track of how many instances of a given complex the system contains. Such an analysis (which involves running a graph isomorphism to establish whether two complexes are the same) is only performed at specific time points when the simulator reports on the contents of the mixture. The memory usage of the simulator is thus bounded by the number of individual agents present, a number which generally remains roughly constant even if the simulation samples a large set of unique complexes.

We briefly sketch the basic approach to computing firing probabilities for rules. The Gillespie procedure is usually applied to interaction networks based on molecular species and the reactions that relate them. In Kappa, however, the conceptual units are patterns along with the rules that relate them. Recall that the agent A on the LHS of rule R in Figure 11 of the main text is

(usually) a pattern, not a molecular species, since **A** might possess many sites other than **s**. The rule simply states that the binding between **A** and **B** occurs independently of whatever else **A** and **B** are bound to. This is precisely the empirical content of the cSIN.

The activity  $\alpha_i$  of a *reaction*  $i$  is a mass action term that is a function of the number of instances in the system of each molecular species appearing on the LHS of reaction  $i$ . In the case of a *rule*  $i$ , we need to count the number of occurrences of the LHS pattern in the large graph that represents the mixture. This amounts to counting the many ways the (usually small) LHS graph of rule  $i$  can be embedded (matched) in the (usually large) mixture graph at time  $t$ . Let this number be  $\theta_i(t)$ . For example, the rule **R** in Figure 11 of the main text will have  $n_{A(s)}n_{B(p)}$  embeddings in a mixture that contains  $n_{A(s)}$  agents of type **A** whose site **s** is free and  $n_{B(p)}$  agents of type **B** whose site **p** is free. Given a stochastic rate constant  $\gamma_i$ , the activity  $\alpha_i$  of rule  $i$  is

$$\alpha_i(t) = \gamma_i \theta_i(t). \tag{1}$$

At the beginning of the simulation we incur an upfront cost by computing all matchings for all rules in the system, which are stored as pointers in appropriate data structures. Once the initial activity of each rule has been calculated, we proceed in the standard fashion:

1. Select a rule  $r$  with probability  $\alpha_r/\lambda$ , where  $\lambda = \sum_r \alpha_r$  is the total system activity.
2. Draw a time advance  $\delta t$  exponentially distributed with parameter  $\lambda$ ,  $p(\delta t) = \lambda \exp(-\lambda \delta t)$ , and advance the simulated wall-clock time  $t$  to  $t + \delta t$ .
3. Select at random one matching of the rule  $r$  chosen in step (1).
4. Execute the action represented by the rule  $r$  selected in step (1).
5. Update the data structure of current matchings to reflect the loss and gain of matchings from the executed action.
6. Update affected rule activities and recompute the overall system activity  $\lambda$ . This is done very efficiently, taking into account causal relationships between rules [11].
7. Repeat.

## 2.2 Time advance with null events

A null event is an attempted rule application that is rejected on the basis of constraints that refer to aspects of an interaction that can only be known at run-time when the rule instantiates as a reaction involving a particular combination of reactants. The point of the Doob-Gillespie method is to avoid null events altogether, since the choice in step (1) of the core loop (section 2.1) is intended to be among reactive events. Various constraints applied at runtime (such as the NR and SR constraints) may reject attempted reactions on the basis of properties that the products would have. Null events always incur a computational cost logarithmic in the number of choices in step (1) of section 2.1 above. But the question is how to advance simulated time when a null event has occurred.

The upshot is that when a null event occurs, simulated time must advance exactly as if the event had been reactive even though no change occurs in the system. This approach has been proposed by Yang et al. [12], but without offering a proof of correctness, which we shall supply here for the sake of thoroughness.

A rule  $i$  is a source of null events if some attempts of applying it will be rejected. The activity of a rule  $i$  at time  $t$  is defined in terms of the set of potential next events—the occurrences of the rule’s LHS in the mixture at time  $t$ ,  $\theta_i(t)$ , see equation 1. We do not know which potential next events would be rejected by the constraint, as knowing them in advance would require computing the constraint for all of them, which is computationally much too expensive—in particular when rejection is rare. Rather, we work with the apparent activity of rule  $i$ ,  $\alpha'_i$ , which is always equal or greater than the true activity  $\alpha_i$  based on potential next events that would be productive:  $\alpha_i = \alpha'_i - \epsilon_i \alpha'_i$ , with  $0 \leq \epsilon_i \leq 1$  being the fraction, or likelihood, of potential null events among the potential next events induced by rule  $i$ :

$$\epsilon_i \equiv \frac{\alpha'_i - \alpha_i}{\alpha'_i}. \quad (2)$$

Let also  $\lambda' \equiv \sum_r \alpha'_r$  be the total apparent activity of the mixture at time  $t$ .

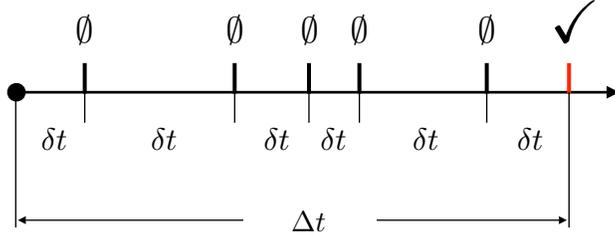
First, if we choose the next rule  $i$  with probability  $\alpha'_i/\lambda'$ , then the probability that rule  $i$  is the next *productive* rule chosen is  $\alpha_i/\lambda$ . This can be seen by considering that, if rule  $i$  is chosen, the probability of choosing a productive copy of that rule is  $\alpha_i/\alpha'_i$  ( $= 1 - \epsilon_i$ ). The probability that rule  $i$  is chosen and productive is thus  $(\alpha'_i/\lambda')(\alpha_i/\alpha'_i) = \alpha_i/\lambda'$ . The probability that the next productive rule is rule  $i$  is just the probability of choosing a productive copy of rule  $i$  divided by the total probability of choosing a productive rule. Since the total probability of choosing a productive rule is  $\sum_i \alpha_i/\lambda' = \lambda/\lambda'$ , we obtain  $(\alpha_i/\lambda')(\lambda'/\lambda) = \alpha_i/\lambda$ . Note that the probability of choosing the next productive reaction to be an application of any rule  $i$  is thus equivalent to the case one would obtain in the absence of clashes.

Next we must show that the distribution  $p(\Delta t)$  of waiting times  $\Delta t$  between *productive* events,  $p(\Delta t) = \lambda \exp(-\lambda \Delta t)$ , results from using the distribution  $p(\delta t)$  of waiting times between any events (productive and null) using the apparent total activity  $\lambda'$ ,  $p(\delta t) = \lambda' \exp(-\lambda' \delta t)$ .

The distribution of waiting times until the next productive event can be written as

$$p(\Delta t) = \sum_{n=0}^{\infty} p(n) \Gamma(\Delta t, n + 1, \lambda'), \quad (3)$$

where  $n$  sums over the number of null events that might occur before a productive event is chosen;  $p(n)$  is the probability of choosing a sequence of  $n$  null events terminated by a productive event. Each of these events generates a time advance  $\delta t$  (Figure 3) based on the same exponential density  $p(\delta t) = \lambda' \exp(-\lambda' \delta t)$ , since the total system activity  $\lambda'$  remains unchanged throughout the null series up and including the productive event (after which  $\lambda'$  may change). The probability density for the sum of  $n + 1$  independent exponentially distributed random variables is



**Figure 3. Null events.** The time line of a sequence of five attempted events, all rejected on the basis of failing some test, and a sixth event that is reactive (i.e. it alters the state of the system). Each increment  $\delta t$  (including the last one leading up to the reactive event) is independently and identically distributed with an exponential probability density  $p(\delta t) = \lambda' \exp(-\lambda' \delta t)$ , as the parameter  $\lambda'$  is unaffected by a “non-event”. The probability density of the sum of all these increments is therefore a Gamma distribution.

a Gamma distribution. In our specific case:

$$\Gamma(\Delta t, n + 1, \lambda') = \lambda' \frac{(\lambda' \Delta t)^n}{n!} e^{-\lambda' \Delta t}. \quad (4)$$

The probability of getting a series of  $n$  null events followed by a productive event is given by  $p(n) = \epsilon^n (1 - \epsilon)$ , where  $\epsilon \equiv \frac{\lambda' - \lambda}{\lambda'}$ . Putting all of this together, we can rewrite equation 3 as:

$$p(\Delta t) = \sum_{n=0}^{\infty} \epsilon^n (1 - \epsilon) \lambda' \frac{(\lambda' \Delta t)^n}{n!} e^{-\lambda' \Delta t} = (1 - \epsilon) \lambda' e^{-\lambda' \Delta t} \sum_{n=0}^{\infty} \frac{(\epsilon \lambda' \Delta t)^n}{n!}. \quad (5)$$

The sum on the RHS of equation 5 is just the expansion for the exponential. This gives

$$p(\Delta t) = (1 - \epsilon) \lambda' e^{-\lambda' \Delta t} \left( e^{\epsilon \lambda' \Delta t} \right) = (1 - \epsilon) \lambda' e^{-(1-\epsilon) \lambda' \Delta t} = \lambda e^{-\lambda \Delta t}. \quad (6)$$

We thus have that the probability of choosing a productive rule  $i$  is  $\alpha_i / \lambda$  and that the waiting time distribution  $p(\Delta t) = \lambda e^{-\lambda \Delta t}$ , demonstrating that the approach to null events taken here results in a case equivalent to the Gillespie-Doob approach for the set of productive reactions.

### 3 Counting complexes in acyclic contact maps

To count the possible species that can be realized from a collection of binding rules with an *acyclic* contact map, we recast the contact map in terms of “views” that are local to each agent. (For a more general analysis and application of the local view concept, see [13].) In essence, creating local views amounts to prying each agent out of the contact map while retaining “bond stubs” (half edges) pointing to the agent and the site at the other end of the bond. Thus, if agent A connects at site  $\mathbf{a}_1$  to agent B at site  $\mathbf{b}_1$  and at site  $\mathbf{a}_2$  to agent C at site  $\mathbf{c}_1$ , we can define a “fragment”  $\mathbf{A}(\mathbf{a}_1^{\mathbf{b}_1 @ \mathbf{B}}, \mathbf{a}_2^{\mathbf{c}_1 @ \mathbf{C}})$  that contains two bond stubs representing surfaces complementary to

fragments  $\mathbb{B}(\mathbf{b}_1^{\mathbf{a}_1^{\text{@A}}}, \dots)$  and  $\mathbb{C}(\mathbf{c}_1^{\mathbf{a}_2^{\text{@A}}}, \dots)$ . For each agent we also generate all versions in which any number of sites lacks a stub, meaning that they will remain unbound. Two complementary surfaces can be plugged together, creating a larger fragment in which these surfaces are unavailable for further interaction. Surfaces that have not been paired with their complement are termed “open”, otherwise they are “closed”. The set of open surfaces constitutes the interface of a fragment. These definitions extend to complexes.

Our counting procedure, specified below, combines fragments  $\mathcal{F}$  into bigger fragments (with new interfaces), starting with the initially provided local views. At every iteration there is a “current set” of fragments that constitutes the material for further combinations in the next iteration. Yet, rather than thinking in terms of fragments, we shall think in terms of interfaces. At the start of the process, we collect all fragments (which at that point are local views) with the same interface into a set and determine its size. This results in a collection of interfaces  $\mathcal{I}$ , each associated with a cardinality  $|\mathcal{I}|$  reporting the number of fragments with that interface. Instead of combining fragments, we combine interfaces and update the associated cardinalities, thus keeping track of the number of fragments with that interface generated up to the current iteration.

The trick is to exploit the fact that any complex is, by assumption, an acyclic complex, which is to say a tree structure. To systematically build a tree, we grow it by repeatedly combining two interfaces, one of which we require to have *exactly one surface*. This can always be done, since an acyclic graph has at least one node with at most one edge. Hence we start with the terminal nodes of the graph (which are connected by exactly one edge or surface). If the other node, call it  $\mathbf{Y}$ , has two surfaces, then this combination will use up one, leaving a fragment with only one surface that can be used to further grow the graph in the next step. If  $\mathbf{Y}$  has multiple surfaces, we build the acyclic subgraphs connecting to each one of these surfaces in exactly the same manner, starting from their terminals. Once each subgraph is completed it has exactly one available surface to connect with  $\mathbf{Y}$ . When all but one surfaces of  $\mathbf{Y}$  are occupied, our tree has grown to include  $\mathbf{Y}$  and has one surface left to be combined with another acyclic subgraph. When all surfaces of a growing tree are exhausted, we have a complex. Since we are not combining fragments but interfaces, we are effectively operating with sets (whose members, however, are never explicitly represented!), enabling us to easily keep track of the cardinality of the interface resulting upon combination. (Do not confuse the cardinality of an interface, i.e. the number of implicitly generated fragments that possess that interface, with the size of an interface, i.e. the number of its surfaces.) Keeping these preliminaries in mind, our procedure reads as follows.

1. From the initial set of local views, construct a list of interfaces, each of which is associated with a cardinality reporting how many local views (fragments) exist with that interface.
2. Loop over all interfaces  $\mathcal{I}$  that have *exactly one surface*. If none are present, go to step 4. If a fragment with interface  $\mathcal{I}$  is to appear in a complex, it must interlock with some fragment with interface  $\mathcal{J}$  containing a surface complementary to  $\mathcal{I}$ 's. Thus, we combine  $\mathcal{I}$  with each interface  $\mathcal{J}$  that has a complementary surface, generating an interface  $\mathcal{K}$ .

The interface  $\mathcal{I}$  is discarded (but not  $\mathcal{J}$ ) and the interface  $\mathcal{K}$  is added to the current set. The cardinality associated with  $\mathcal{K}$  is updated to  $|\mathcal{I}| \cdot |\mathcal{J}| + |\mathcal{K}|$ , where  $|\mathcal{K}|$  is the prior cardinality of  $\mathcal{K}$ , if  $\mathcal{K}$  already existed in the prior set of interfaces.

3. Repeat step 2.

4. The current set contains an empty interface  $\varepsilon$ , and its cardinality is the number of possible complexes that can be formed with the local views obtained from the initial contact map.

By using arbitrary-precision arithmetic, we can *exactly* determine the astronomically large numbers of complexes that arise from the artificial contact maps generated in the next section.

We have omitted one final issue: symmetries (automorphisms). The only symmetry that can arise in an acyclic complex stems from dimerization, that is, the binding of two copies of the same agent type at sites with the same label, as in the dimerization of many receptors in signaling cascades. We refer to this as “self-binding”, because it shows up as a loop in the contact map. (Notice that a loop returning to the same site at which it originated is not a cycle, as it cannot be realized as a cyclical complex; rather, such a loop results in a dimer.) This case is easily taken care of by redefining “interfaces of size 1” in step 2 of our procedure as interfaces that have one non-self-binding surface and possibly one self-binding surface. The combinations on the non-self-binding surface proceed as described above. In addition, we need to account for combinations on the self-binding surface, and these contribute  $|\mathcal{I}| \cdot (|\mathcal{I}| - 1)/2$  to the cardinality of the resulting interface.

## 4 Random Acyclic Graphs

We generate Random Acyclic Graphs (RAGs) using the following iterative procedure. One begins at step 1 with a pair of nodes (named  $A_0$  and  $A_1$ ), each with a single site labeled  $A_0.s_1$  and  $A_1.s_1$ , respectively, connected by an edge (representing a possible binding action). At each subsequent step  $n$ , a new node  $A_n$  is added to the graph.  $A_n$  is created with a single site,  $A_n.s_n$ . A node  $A_i$  is then chosen at random with equal probability from the set of existing nodes  $\{A_0, \dots, A_{n-1}\}$ , and a new site is created on  $A_i$ , labelled  $A_i.s_n$ . Finally, an edge is created between  $A_i.s_n$  and  $A_n.s_n$ . By construction, this procedure results in an acyclic site graph. A so-generated graph with  $N$  nodes has  $N - 1$  edges, which corresponds to a fairly low edge density  $\rho = 2N^{-1}$ . Indeed, with  $N = 167$  nodes  $\rho \approx 0.012$ , which is considerably smaller than the edge density observed for the cSIN ( $\rho \approx 0.039$ ).

To provide a more accurate comparison to the cSIN, we extend the previous procedure by a second phase in which more edges are added. In phase 1, a graph is generated with  $N$  nodes, exactly as previously described. In phase 2, new edges are added to this graph in a manner that avoids proper cycles (recall that a proper cycle in a site graph is a path that starts and ends at the same node and consists of bonds that do not touch the same site twice, section “Methods”  $\rightarrow$  “Dealing with network cycles” in the main text). To set up phase 2, define  $\mathcal{S}_n$  to be the set of all sites that existed prior to step  $n$  in phase 1,  $\mathcal{S}_n = \{A_j.s_k \mid 0 \leq j < n, 1 \leq k < n\}$  (with the proviso that  $A_j.s_k$  exists, since the sites of  $A_j$  are not labelled consecutively). Also, define  $\mathcal{S}_{A_i}$  as the set of sites belonging to a particular node  $A_i$ . To add new edges, we first choose a node  $A_i$  at random (uniformly) from the set  $\{A_2, \dots, A_N\}$ . We then choose another node  $A_j$  from the set  $\{A_2, \dots, A_{i-1}\}$  (i.e. a node *older* than  $A_i$ ). Finally, we choose a site  $s_k$  at random with uniform probability from the set  $\mathcal{S}_i \cap \mathcal{S}_{A_j}$ , that is, a site  $A_j.s_k$  from the sites on  $A_j$  that existed prior to step  $i$ . An edge is then placed between  $A_i.s_i$  (the only site on  $A_i$  that exists at step  $i$  of phase 1) and  $A_j.s_k$ ; if this edge already exists, another node is chosen at random from the set

$\{A_2, \dots, A_{i-1}\}$ , until a new edge is successfully placed or all possibilities in  $\mathcal{S}_i$  have been exhausted; at that point a different  $A_i$  is chosen at random. These steps are repeated until the edge density of the RAG is approximately equal to the desired edge density of the cSIN,  $\rho \approx 0.039$ .

This way of adding edges cannot produce contact maps with proper cycles. Any newly added edge in phase 2 is always between a site  $A_i.s_i$  ( $2 \leq i \leq N$ ) and a node  $A_j$  with  $j < i$  that belongs, by construction, to an acyclic subgraph generated before step  $i$  in phase 1. The only “return” from a node in that subgraph back to  $A_i$  was created exactly at step  $i$  and involves the same site  $A_i.s_i$ , and is, therefore, not a proper cycle. A complex constructed from such a contact map can never be a ring, since in a complex (unlike in a contact map) any site can be bound at most once (in particular  $A_i.s_i$ ), see definitions in section 1.2.

In practice we find that this procedure efficiently generates RAGs for  $N \gtrsim 60$  at an edge density similar to that of the cSIN. We create sets of RAGs with varying number  $N$  of nodes but a fixed edge density  $\rho \approx 0.039$ . Each point in Figure 4B of the main text reports the average number of possible complexes as determined from 10 independently generated RAGs with a given  $N$ .

## 5 The effect of size constraints on the number of possible complexes

In the above analysis, the assumption is made that all possible complexes can be physically realized by the proteins and their interfaces represented in the (acyclic) cSIN-like contact maps. Although we include steric constraints at runtime (particularly in the “stable rings” scenario, section 6), the counting algorithm described in section 3 does not account for steric effects that might prevent the formation of certain complexes. In this section we perform a simple calculation to explore the consequences that steric constraints might have on the total number of molecular species that an interaction network could form. The case we consider here represents a fairly strong constraint, in which steric effects become more and more prominent as complexes get larger. Given that the surface area of a complex will tend to increase with increasing size, this might not represent the most realistic situation. We nonetheless consider this model to demonstrate that even strong steric constraints do not curtail combinatorial complexity significantly, unless the parameters of the model are set to extreme values.

One approach to assessing the reduction in the number of complexes from higher-order steric constraints would be to enumerate the set of possible complexes and then remove some fraction of these (as a function of complex size). This fraction would represent a parameter of the model. Such an analysis cannot be performed directly, because the large number of possible complexes prevents their explicit enumeration. However, based on the initial explicit enumeration depicted in Figure 4A of the main text, we argue that the total number of complexes scales (approximately) in an exponential fashion with complex size, leading us to posit:

$$N = \sum_{s=1}^M N(s) = \sum_{s=1}^M be^{as}, \tag{7}$$

where  $N$  is the total number of complexes that can be formed by the network,  $N(s)$  is the number of complexes of size  $s$  with  $s$  ranging from 1 (monomers) to  $M$ , the size of the largest complex that the network could form.  $a$  and  $b$  are free parameters. Computing the geometric

series in equation 7, we can write:

$$N = b \frac{e^{aM} - e^a}{e^a - 1}. \tag{8}$$

The model represented by equation 7 states that  $N(1) = be^a$  and  $N(s) = e^a N(s - 1)$ , expressing an exponential progression as one forms complexes of size  $s + 1$  from complexes of size  $s$ . Let us now assume that the addition of a protein in growing complexes from size  $s$  to  $s + 1$  encounters a steric hindrance that reduces the number of attainable complexes by a factor of  $p$  ( $0 < p \leq 1$ ). This assumption means that  $N(1) = be^a$  and  $N(s) = pe^a N(s - 1) = bp^s e^{as}$ . Thus,

$$N = \sum_{s=1}^M bp^s e^{as} = b \frac{e^{(a+\log p)M} - e^{a+\log p}}{e^{a+\log p} - 1} \tag{9}$$

where we see that the effect of  $p$  is to reduce the effective value of  $a$  by  $|\log p|$ .

To provide a sense for how large this effect might be, we estimate values of  $b$  and  $a$  for the RAGs discussed in section 4. Taking a network with 167 nodes, we know that the number of monomers will be 167; this gives us  $b = 167e^{-a}$ . For RAGs without any steric constraints (i.e. with  $p = 1$ ), we have that  $N \approx 10^{40}$  (see Figure 2 of the main text). Assuming that  $M \approx 40$  (see Figure 6 of the main text) allows us to solve equation 8 numerically for  $a$ . Using Mathematica [14], we get  $a \approx 2.23$  and thus  $b \approx 17.96$ . Using these values, Figure 4 depicts how the total number of complexes will change as  $p$  is decreased (thus increasing the steric hindrance at every growth step).

From Figure 4 we can see that even when  $p$  is fairly small the network is still likely to produce a very large number of possible complexes. Even if only 20% of complexes of a given size can be realized, the total number is still  $\sim 10^{12}$ . This approximate calculation suggests that steric constraints would have to be incredibly strong in order to reduce the number of molecular possibilities to numbers that would allow their simultaneous sampling by a cell.

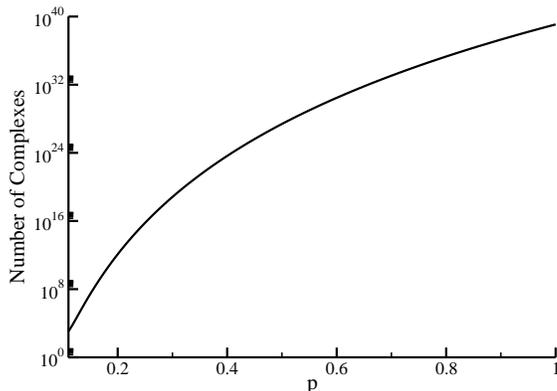
## 6 The “stable rings” scenario

### 6.1 The thermodynamic rationale for the “stable rings” scenario

The SR scenario posits that all cycles in the contact map of the cSIN correspond to physical structures whose bonds along the cycle backbone can be satisfied simultaneously whenever a complex containing the elements of the cycle arises, regardless of cycle length. In order to implement the SR case, we must first calculate the rate at which cycles should close (i.e. the rate at which the intramolecular A-B bond missing on the left of Figure 11B in the main text should form). We follow the work of Saiz and Vilar [15] by defining the standard free energy change upon binding of two proteins ( $\Delta G_b^0$ ) as:

$$\Delta G_b^0 = \Delta G_p^0 + \Delta G_i^0, \tag{10}$$

where  $\Delta G_p^0$  represents the positional entropy loss entailed when taking two proteins that can freely diffuse around a particular molar volume and confining them to a given binary complex.



**Figure 4. Effect of steric constraints.** The plot shows the effects that higher-order geometric and steric constraints might have on the total number of possible complexes a network could form. In this case, we assume that the number of complexes increases exponentially with complex size up to some maximal complex size.  $p$  represents the fraction of complexes that cannot be formed due to steric constraints when attempting to bind a protein to complexes one size smaller. The black line is calculated using equation 9, with parameters  $a = 2.23$ ,  $b = 17.96$  and  $M = 40$  chosen to approximate the case of a RAG with  $\sim 167$  nodes. Note that a sizable fraction ( $\sim 89\%$ ) of complexes must be sterically prevented *at each size* in order to produce numbers similar to the total number of unique complexes observed in a single simulation (i.e.  $10^4$ , Figure 5A in the main text).

$\Delta G_i^0$  represents the free energy of the specific molecular interactions in the complex, including contributions from the desolvation of the two protein interfaces and the molecular contacts (e.g. electrostatic and Van der Waals interactions) formed upon binding. We assume that  $\Delta G_b^0 < 0$  for every pair of interacting surfaces in the cSIN *independently*; that is, for every edge in the graph we assume that the implied binding reaction will favor the bound form even when no other proteins from the cSIN are included in the system. As discussed in “Methods” (main text), we employ various interaction affinities, but here we focus on a case in which every reaction has a dissociation constant of 10 nM. If we assume the reactions are taking place at an absolute temperature of  $\sim 300$  K, and that the standard molar positional entropy loss upon binding is  $\sim 9$  kcal mol $^{-1}$  [16], a 10 nM dissociation constant implies that  $\Delta G_i^0 \sim -21$  kcal mol $^{-1}$ .

The above calculation is based on the assumption that two proteins are binding to one another in a bimolecular fashion. Cycle closure, however, is a unimolecular reaction. Moreover, in the SR case, we assume that complexes can structurally satisfy all of the interactions in the cycle simultaneously. This implies that the positional entropy loss of the cycle closure reaction is  $\sim 0$ ; that is, all the members of a particular complex are already constrained by being bound to one another, such that the formation of the final bond does little to change the positional entropy of either protein that participates in it. The change in free energy represented by this unary reaction is thus  $\Delta G_b^0 \sim \Delta G_i^0 \sim -21$  kcal mol $^{-1}$ . This implies that such rings are very stable, as has been

argued elsewhere [15, 17]. Indeed, if we designate the forward rate of cycle closure as  $u_+$  and the cycle opening rate as  $u_-$ , then the free energy of cycle closure implies that  $u_-/u_+ \sim 10^{-16}$ .

If we assume that the cycle opening rate and the off rate  $k_-$  of the binary reaction are approximately equal, we have  $u_-/k_+ = 10^{-8}$ , so the binary binding rate  $k_+$  must be much smaller than the cycle closure rate:  $u_+ \sim 10^8 k_+$ . Our reasoning, so far, applies to deterministic rate constants. In a stochastic setting, bimolecular rate constants, such as  $k_+$ , are volume dependent [10]. Let  $\kappa_+$  denote the stochastic binding rate constant expressed in units of  $\text{molecule}^{-1}\text{s}^{-1}$ . If the deterministic  $k_+$  is in the usual units of  $\text{M}^{-1}\text{s}^{-1}$ , then  $\kappa_+ = k_+/(N_A V)$  with  $N_A$  denoting Avogadro’s number and  $V$  the reaction volume. In the present work, we take  $V = 4.2 \cdot 10^{-14}$  L as the volume of a haploid yeast cell [18], which yields  $u_+ \sim 10^{18} \kappa_+$ .

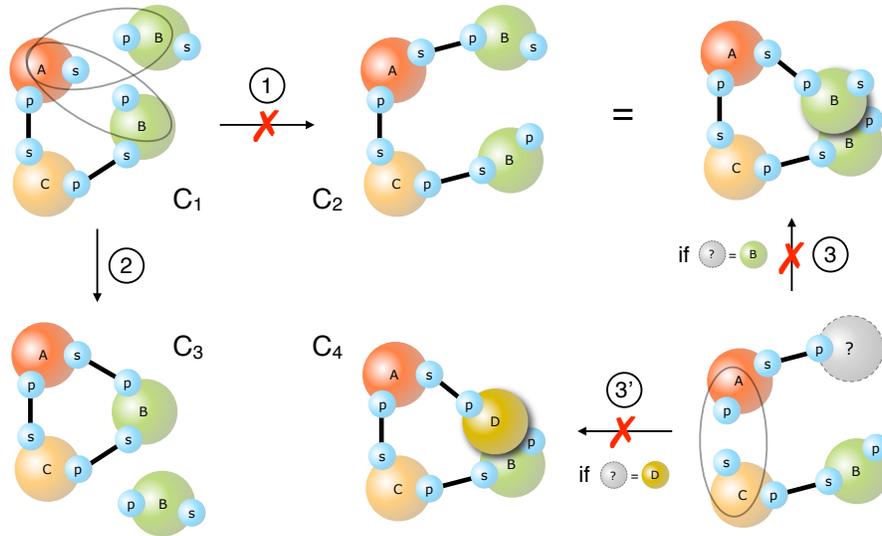
We are led to conclude that cycle closure is so fast that, were it to be included explicitly in the simulation, most of the events at steady state would consist of cycle openings and closings. Moreover, once formed, cycles are so stable that they will tend to persist. Using the parameter values of the main text (“Results” → “Network Dynamics”), the time scale of the full excision of one member from its cycle (“cycle decay”) is likely to be on the order of  $10^{14}$  time units. Our longest simulations stretch over  $\sim 100$  time units, meaning that on the time scale relevant to this study, a cyclic structure will not decay once it has formed. This is true not only for the 10 nM dissociation constant used in our illustrative calculation, but for every interaction affinity case discussed in the main text. We therefore implement cycle closure at an infinite rate as in the work of [17], i.e. an SR simulation is defined by  $u_-/u_+ = 0$ .

One intriguing consequence of the thermodynamics of cycles is the fact that a set of proteins can form a very stable ring even if one or more of the bonds in the cycle is fairly weak. This leads to a situation in which a high-throughput screen, such as Yeast-2-Hybrid, might miss interactions that are not stable on their own but are stable in the context of a proper cycle. We discuss this possibility further in section 8.5.

## 6.2 The implementation of the “stable rings” scenario

An implementation of the SR policy at the level of local Kappa rules would require an enumeration of all possible cycle closures. Even if such rules could be generated automatically, their number would be huge. We therefore implement the SR policy as a global directive to the simulator. Whenever a binding reaction occurs, the newly formed complex is explored to determine if it contains free sites that could bind any free sites on the protein that just entered the complex. If one or more such sites exist, they bind one another instantly without advancing simulated wall-clock time, which corresponds to a reaction with an infinite forward rate [10]. If a site (or sites) can participate in more than one cycle within the complex under consideration, the cycles that are actually formed are chosen at random with equal probability. Once a cycle is formed, any bonds in it are prevented from matching a dissociation rule, since such an event would immediately be reversed by an infinite-rate cycle closure event.

Setting the rate of cycle closure to infinity prevents many polymerization reactions from occurring, since complexes like the open triangle  $C_1$  in Figure 5 immediately react along route 2 and never linger to pick up another B along route 1 to form  $C_2$ , which nucleates a polymer. Yet,



**Figure 5. The excluded volume constraint.** The left of the figure illustrates the infinite cycle closure condition.  $C_1$  is a complex in which protein A can bind protein B. Because of the local nature of binding rules (corresponding to the local nature of the information in the SIN), there are two possibilities, here indicated by ovals: A binds the B that is already in the complex (unimolecular reaction 2), or A binds a B that is not in the complex (bimolecular reaction 1). Based on thermodynamic arguments (section 6.1), the SR constraint posits that whenever a complex contains members that could form a cycle (as we assume to be the case for protein agents A, B, and C), the cycle closure occurs immediately (reaction 2), thus preventing reaction 1. The undesirable complex  $C_2$  is also an intermediate towards polymerization. (In the case of an ambiguity in the molecularity of a local binding interaction, it is always the case that the unimolecular reaction is a ring closure, whereas the bimolecular interaction may be a step towards polymerization, depending on the context.) However, the polymeric intermediate  $C_2$  could also be formed by reaction 3, which involves the binding (indicated by the oval) of two dimers that would bring together into one complex all the members of a cycle, yet with one of the sites required for cycle closure (here  $s$  of A) already occupied. This is a steric contradiction that must be prevented, regardless of whether the outcome is a polymerization intermediate (as when A is bound to a B), reaction 3, or not (as when A is bound to a D), reaction 3'. Although not a polymerization intermediate, the latter product is prevented nonetheless because of steric consistency. The SR scenario is therefore more stringent than the sole prevention of polymers.

the problem is not fully eliminated, as there is another route to  $C_2$ . A local binding interaction between the dimers  $A(p, s^1), B(p^1, s)$  and  $C(s, p^1), B(s^1, p)$  on the RHS of reaction 3 in Figure 5 also generates the troublesome  $C_2$ . This reaction must be prevented on the basis of a steric consistency argument. Recall that the assumption of the SR scenario states that whenever all members of a cycle are in the same complex their geometry satisfies all their bindings. Hence a dimerization as in reaction 3 of Figure 5 would be sterically inconsistent, as it brings together all members of a possible cycle (proteins A, B, and C) with one member (A) already being occupied at a binding site ( $s$ ) that must, by assumption, be in steric juxtaposition with its binding partner in

the potential cycle (here the B bound to C). Consistency, therefore, requires that the B bound to A (or the B bound to C) prevents reaction 3, as it would end up occupying the same volume already occupied by the other B. This is suggested in Figure 5 by a pictorial clash in the rendering of  $C_2$  emerging from reaction 3, but the reader should not forget that these diagrams are graphs with no geometric content. The same argument holds, of course, if A were bound at site  $\mathbf{s}$  to a protein other than B, say D, as in reaction 3' of Figure 5. The offense to the SR assumption comes from the occupancy of any site (within an attempted complex) that *would have* bound at infinite rate another site in the complex to close a cycle. We refer to this constraint as the “excluded volume” test. It, too, is implemented as a directive to the simulator: Any attempted binding is inspected at run-time to determine whether the resulting complex contains sites that offend the SR assumption in the sense just described. Simulations denoted as “stable ring” scenarios in this supplement are performed using the set of binary local rules derived from the cSIN (e.g. rules as in equation 2 of the main text) together with two amendments applied at run-time: (i) cycles close infinitely fast and are forever stable, and (ii) attempted complexes must pass the excluded volume test. Rejections of attempted reactions based on (ii) are treated as null events, see section 2.2.

Using Figure 5 as a guide, the infinite closure of cycles means that as soon as B binds C to form  $C_1$ ,  $C_2$  occurs.  $C_1$  never enters the mixture. Structurally, this means that the incoming B really forms two contacts at once (by our idealized assumption of perfect juxtaposition). However, based on the binding rules in the system, there are two independent ways for B to close the ring: by first binding to C (as on the left of Figure 5) followed by immediate ring closure or by first binding to A followed by immediate ring closure. The independence adds up to twice as fast a ring “discovery” kinetics than is warranted by our assumption that both bonds form at the same time. To recapture the correct kinetics, each time an agent binds a complex in an event that happens to create the opportunity for a ring closure, we reject the proposed binding of the agent with probability 1/2. In this way, we effectively divide the ring discovery rate by 2 without the need of knowing a priori which events are potential ring closures, which is combinatorially prohibitive. The price, again, is the creation of null events, see section 2.2.

A consequence of preventing polymerization in the SR scenario is the avoidance of any molecularity ambiguity of local binary rules. In the SR scenario, productive reactions induced by binary rules will always be bimolecular, and are therefore assigned a single bimolecular stochastic rate constant.

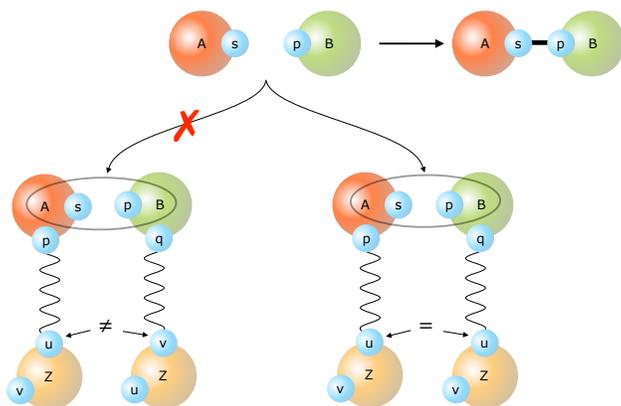
## 7 The “no rings” scenario

The rationale behind this scenario is to *be agnostic about the nature of steric constraints that prevent the formation of polymers*. We simply prevent polymers by fiat, using a non-local constraint illustrated in Figure 6.

### 7.1 The implementation of the “no rings” scenario

Polymerization requires the coming together of parts, each of which with at least one copy of the same protein as the other. But the reverse is not true: it is not the case that the binding of parts that share a protein type necessarily leads to polymerization. It depends on how the the copies of

the common type are bound inside the complex. Figure 6 depicts a binary local rule representing the binding of A and B. Suppose the LHS pattern of the rule is matched (grey oval) in the mixture by a pair of complexes as shown on the lower left of Figure 6, both of which contain a copy of protein Z (a protein with at least two binding sites). The wiggly lines stand for an arbitrary intervening set of connected proteins in each complex. The complex resulting from the binding of



**Figure 6. The no-polymerization constraint.** The figure illustrates a syntactical criterion for rejecting binding interactions that would lead to the formation of polymers: The binding of complexes with duplicate parts (here Z) must be prevented whenever there is a path connecting them on *different* sites. See the text for a detailed explanation.

these two parts obviously contains two Z's for which there is a connecting path – a path being a sequence of bonds representing an instruction for traveling through the complex from a particular site at a source protein to a particular site at a target protein. If a connecting path between the two copies of protein Z starts and ends at different sites of Z, the complex constitutes an intermediate in a polymerization process and must be prevented by rejecting the attempted binding. To see this, assume, as shown in Figure 6, that Z is connected by some path to A starting at Z's site u, in pseudo-notation:  $v.Z.u \rightsquigarrow A$ , while u is free on the copy of Z in the other part,  $B \rightsquigarrow v.Z.u$ . Clearly, the new complex  $v.Z.u \rightsquigarrow A-B \rightsquigarrow v.Z.u$  could grow another copy of  $\rightsquigarrow A$  on Z's free u site. The very existence of the  $v.Z.u \rightsquigarrow A$  complex proves that this is possible given the set of local rules available to the system. Yet, if both copies of Z are tied up at the same site u, as shown on the lower right of Figure 6, they could not function as a platform for growing a polymer.

As a consequence of this constraint, complex  $C_4$  in Figure 5 is allowed, but not  $C_2$ . However, an “isomer” of  $C_2$ , in which A is bound to B on its site s passes the no-polymerization test. The constraint depicted in Figure 6 prevents all polymerizing complexes and is more subtle than a simple ban of complexes that contain duplicates of proteins.

Rejected attempts at reaction based on the NR (or SR) constraint result in null events that are handled as detailed in section 2.2. In the NR and SR scenarios we have that binary binding rules always correspond to bimolecular reactions in the mixture, avoiding any potential mismatches in reaction arity (Figure 3 of the main text).

In sum, the NR scenario prevents polymerization by fiat, catching all higher-order geometric constraints that we cannot directly express in Kappa. However, neither the SR case nor the NR case is likely to represent the reality of complex formation in the cell. Some of the cycles in the contact map of the cSIN might represent SR complexes, others might follow the NR scenario or perhaps even give rise to polymers of limited size.

## 7.2 The relationship between “stable rings” and the “no rings” scenario

As noted in section 1.3, polymerization and ring closure are flip sides of the same coin. The trick in the SR scenario is to prevent polymers by invoking the steric and kinetic consequences that arise when proper cycles in the contact map are interpreted geometrically as rings. In the NR scenario, however, we prevent polymers purely logically—through the constraint depicted in Figure 6. Consider that constraint in the context of Figure 5. The no-polymerization criterion prevents the bimolecular reaction 1 between A and site p of B because there already is a path from A to a different site of another B. Note, however, that the same *path pattern* arises upon cycle closure, reaction 2, except that the path starting at B loops back to the same rather than a different instance of B. From a purely formal standpoint, if the constraint is to forbid certain paths, then, perhaps, cycle closure should be prevented as well, as it relies on the same forbidden path pattern. This “formal” argument actually blends with a structural argument. If the juxtaposition between A and B in  $C_1$  (Figure 5) were perfect for ring closure, we would be back in the SR scenario. For the NR case to be different, we are forced to assume that the juxtaposition between A and B is *not* conducive to ring closure. This represents a case in which either the complex is very flexible, making ring closure entropically costly, or the B-binding site of A points in a different direction than the location of the resident B. The latter is compatible with the no-polymerization constraint, which allows A to indeed pick up an additional B (with the proviso on binding sites), see Figure 6 (where Z plays the role of B). Thus, in addition to preventing polymerization, the NR scenario also sets  $u_+ = 0$ , which is implemented by rejecting every binding reaction that would result in the closure of a cycle.

## 8 Additional Results

### 8.1 Alternative definitions of distance

In the main text, we define the “normalized distance” between two cells as:

$$d(i, j) = \frac{|C_i \Delta C_j|}{|C_i \cup C_j|} \tag{11}$$

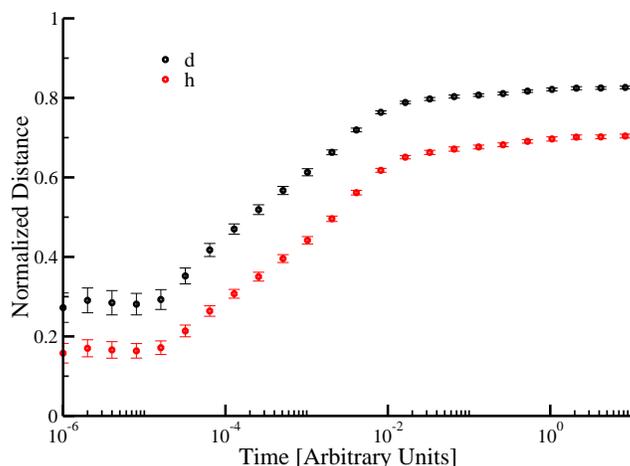
where  $C_i$  is the set of unique complexes in cell  $i$ ,  $C_i \Delta C_j$  is the symmetric difference between  $C_i$  and  $C_j$  (i.e.  $(C_i - C_j) \cup (C_j - C_i)$ ) and  $C_i \cup C_j$  is the union of  $C_i$  and  $C_j$ . This represents a very natural distance since  $d(i, j)$  is just the probability that a complex in either  $i$  or  $j$  will only be found in one of the two cells.

One could consider several alternative definitions of distance. For instance, if we think of the “complexome” of a cell as a string of 1’s and 0’s – with a 1 in a particular sequence position

indicating that the presence of the corresponding complex, and a 0 indicating its absence. In this representation,  $|C_i \Delta C_j|$  is just the ‘‘Hamming Distance’’ between the two complexome strings. It is natural to normalize the Hamming Distance by the distance we would obtain if the sequences were ‘‘orthogonal’’ in the complexes they contain (i.e. if every 1 in cell  $i$  corresponded to a 0 in cell  $j$  and vice versa), which yields:

$$h(i, j) = \frac{|C_i \Delta C_j|}{|C_i| + |C_j|}. \quad (12)$$

This definition is very similar to equation 11 and yields essentially the same general result (see Figure 7.)



**Figure 7. Comparing different definitions of the distance between cells.** This figure is based on NR simulations with a uniform  $K_D$  of 10 nM. Each point represents an average over all unique comparisons between 15 independent simulations and the error bars correspond to  $\approx 95\%$  confidence intervals. The black curve is calculated using equation 11; this data is equivalent to that displayed in Figure 5B of the main text. The red curve is calculated using equation 12. As one would expect, the two curves are very similar, with the normalized Hamming Distance giving distances somewhat smaller than the normalized distance discussed in the text. This derives from the fact that the normalizing factor in the Hamming Distance case is always greater than or equal to the normalization factor for  $d(i, j)$ ; i.e.  $|C_i| + |C_j| \geq |C_i \cup C_j|$ .

Both of the above distances do not, however, weigh differences in the *copy number* of a particular complex. A complex that occurs in both cells  $i$  and  $j$  contributes the same weight to the overlap, regardless of the difference in the number of copies with which it occurs in each of these cells. We assess differences between cells that arise from copy number variations by considering two additional definitions of the distance. In the first case, we define a distance  $H_N$  which is similar to

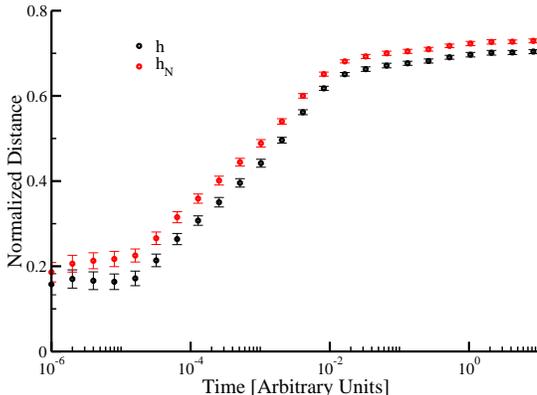
the Hamming distance but contains information about the differences in abundance of complexes:

$$H_N(i, j) = \sum_k \frac{|N_k(i) - N_k(j)|}{N_k(i) + N_k(j)} \tag{13}$$

where  $N_k(x)$  indicates the copy number of a unique complex of type “ $k$ ” in cell “ $x$ ”,  $|N_k(i) - N_k(j)|$  is the absolute value of the difference in copy number for a complex of type  $k$  between cells  $i$  and  $j$ , and  $k$  ranges over all the types of unique complexes in cells  $i$  and  $j$ . Like the Hamming distance, equation 12, each term of the sum in equation 13 assumes a value between 0 and 1. However, two cells that both contain a particular type of complex  $k$  do not contribute to the Hamming distance  $h(i, j)$ , equation 12, whereas such cases will increase the distance  $H_N(i, j)$ , equation 13, whenever  $N_k(i) \neq N_k(j)$ . As in the case of  $h(i, j)$ , two orthogonal cells will have a maximal distance of  $H_N(i, j) = |C_i| + |C_j|$ , which we use as a normalizing factor, yielding:

$$h_N(i, j) = \frac{H_N(i, j)}{|C_i| + |C_j|} \tag{14}$$

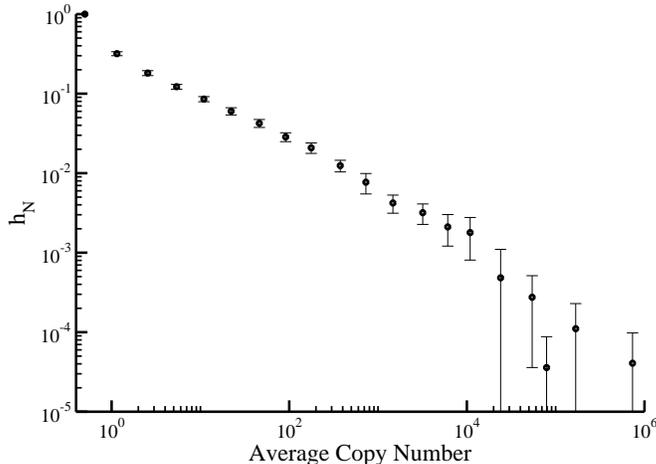
which is analogous to the normalized Hamming distance, equation 12. Equation 14 produces similar results to equation 12 (see Figure 8), with  $h_N > h$  as one would expect.



**Figure 8. Comparing the Hamming distance with a similar distance definition that includes differences in copy number.** This figure is based on NR simulations with a uniform  $K_D$  of 10 nM. Each point represents an average over all unique comparisons between 15 independent simulations and the error bars correspond to  $\approx 95\%$  confidence intervals. The black curve is calculated according to equation 12 (corresponding to “ $h$ ”) and the red curve is calculated using equations 13 and 14 (corresponding to “ $h_N$ ”). The two distances are very similar, with differences in copy number (which are ignored by the definition of  $h$ ) leading to slightly larger distances  $h_N$ .

Given that large complexes tend to be comparatively rare (see Figure 6 in the main text) and distinct between cells (see Figure 5C in the main text), it is useful to stratify complexes into

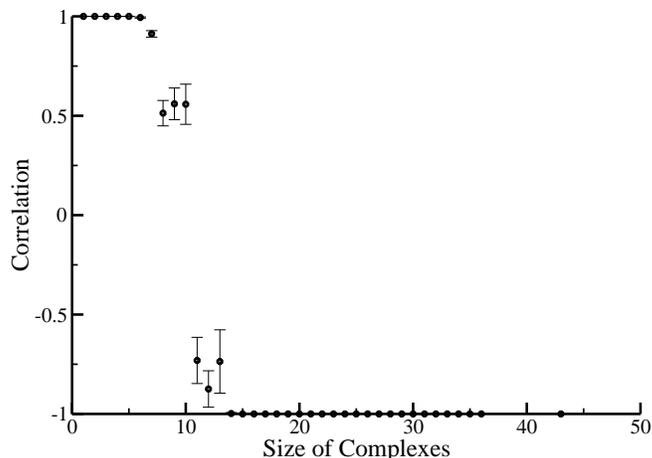
abundance classes and calculate how the distance between a pair of cells varies as a function of abundance class. For two cells  $i$  and  $j$  at time  $t$ , and each complex of type  $k$  in either cell, we obtained the arithmetic average of its copy number as  $\overline{N}_k(i, j) = \frac{1}{2}(N_k(i) + N_k(j))$ . We then binned complexes according to their average copy number, using exponentially distributed bins (e.g.  $0 < \overline{N}_k(i, j) \leq 1$ ,  $1 < \overline{N}_k(i, j) \leq 2$ ,  $2 < \overline{N}_k(i, j) \leq 4$ , etc.) in order to account for the log-normal distribution of copy numbers present in the simulations. Within a given bin, we calculated the average value of  $h_N$  based on the complexes within that bin and plotted the resulting values against the average abundance in that bin, Figure 9. The similarity between cells increases dramatically the higher the abundance class over which the distance  $h_N$  is computed.



**Figure 9. The dependence of  $h_N$ , equation 14, on the copy number of complexes.** This figure is based on NR simulations with a uniform  $K_D$  of 10 nM and is computed at the last time point in Figure 5B in the main text. We calculated the values of  $h_N$  by first binning complexes according to their average copy number in two cells. The first bin, from 0 to 1, contains all of those complexes that have only one copy in one cell in a particular comparison. The bins are sized exponentially, with the upper endpoint of each bin separated from the upper endpoint of the preceding bin by a factor of 2. For each bin, we calculated the value of  $h_N$  for the complexes included in that bin. Each point represents an average over all unique comparisons between 15 independent simulations at the final time point of Figure 5B in the main text. The error bars correspond to  $\approx 95\%$  confidence intervals. As copy numbers increase, the value of  $h_N$  decreases dramatically. The average distance varies approximately as a power law with an exponent of  $-0.7$ .

One can also consider the impact of copy numbers on the similarity between cells by computing the correlation between two cells as a function of complex size. The correlation is calculated in the natural way, by comparing the number of copies of any given complex in cell  $i$  to the number of copies of that same complex in cell  $j$ . As shown in Figure 10, small complexes exhibit a correlation very close to 1, indicating that the copy number of a small complex in one cell is a good predictor of the copy number of that complex in other cells of a population. This is mostly

due to the fact that certain small complexes occur at very high numbers (see Figure 6 in the main text); these very common complexes tend to dominate the correlation. At intermediate sizes, some variation in copy number is observed, while for large complexes (i.e. complexes larger than 10), we obtain correlations  $\sim -1$ . This is due to the fact that cells do not exhibit much overlap in their repertoire of large complexes: if a complex is present in cell  $i$ , it is typically absent in  $j$ , and vice versa, leading to a correlation of  $-1$ .



**Figure 10. The correlation between two cells as a function of complex size.** Each point in this plot represents an average over all unique pairwise comparisons between 15 independent simulations and the error bars correspond to  $\approx 95\%$  confidence intervals. In this case, the data is taken from NR simulations at a  $K_D$  of 10 nM; the comparisons are made between simulations at the final time point in Figure 5A of the main text. The correlation at each complex size is calculated by comparing the copy number of any given complex in a cell  $i$  to the copy number of that same complex in cell  $j$  and computing Pearson’s correlation in the usual fashion. We observe correlations  $\sim 1$  for small complexes while, for large complexes, we find a correlation of  $\sim -1$ . The perfect negative correlation at large complex sizes is due to the fact that, if a large complex is present in cell  $i$ , it will not be present in cell  $j$ , and vice versa.

Taken together, the above results indicate that the difference between cells depends strongly on the size of the complexes over which distance is computed. Small complexes, which tend to exist at high copy numbers, do not vary significantly between cells in terms of their presence or absence (see Figure 5C in the main text), but do exhibit some differences in copy number (Figures 9 and 10). Complexes of intermediate size tend to have considerably smaller copy numbers and somewhat larger variation (either in presence/absence or in copy number) between cells. Individual large complexes tend to be found in very small copy numbers (on average just 1 copy per cell) and each large complex tends to be unique to a given cell. All other alternative definitions of distance that we have examined (including the simple Euclidean distance between the copy number vectors of two cells, the angle between those vectors, etc.) produce exactly the

same results as those discussed here (data not shown).

## 8.2 Structure-based affinities

In the main text we consider three protocols for assigning affinities to the interactions listed in the cSIN. One protocol is based on estimating affinities from available crystal structures of domain-domain interactions. To identify such interactions and their corresponding structures, we must construct *de novo* the cSIN, which we originally extracted from the SIN. We shall refer to the new version as the cSIN2.

The construction of the yeast SIN by Gerstein and coworkers [19] begins with a set of high-confidence interactions defined at the level of whole proteins (i.e. without information regarding sites). For each protein in this set, a domain structure is obtained based on the Pfam database of protein domains [20]. If two proteins interact according to the high-confidence list, the SIN construction checks whether any of their domains interact according to iPfam [21], a domain-domain interaction database derived from an analysis of solved protein structures. If a high-confidence interaction can be reconciled to a domain-domain interaction, the domains are added to the corresponding proteins in the SIN and edges are placed between them. These domains thus represent the “sites” of the SIN (and hence the cSIN), see Figure 2 of the main text.

We estimated the interaction affinity between any two domains in the cSIN on the basis of the protein structures used to define that particular domain-domain interaction in iPfam. To do this we first re-created the cSIN using an updated release of Pfam (we employed Pfam release 21, since this release contained the most recent version of iPfam at the time of this writing). We shall refer to this updated cSIN as cSIN2. We first determine the Pfam domains present in each protein of the original cSIN. We then determine for every interacting pair of proteins their corresponding set of interacting domains, as defined in iPfam. For example, suppose that proteins A and B interact in the cSIN and that domains  $\{k, l, m\}$  are found in A and domains  $\{m, n, o\}$  are found in B. If iPfam asserts that domains of type “l” are known to interact with domains of type “m”, we add a site corresponding to domain l to protein A and a site corresponding to domain m to protein B and draw an edge between these domain instances in the cSIN2. Since an edge is placed between *all* sites that can interact, two proteins may have more than one edge connecting them. The cSIN2 obtained with this updated release of Pfam is very similar to the original cSIN (data not shown).

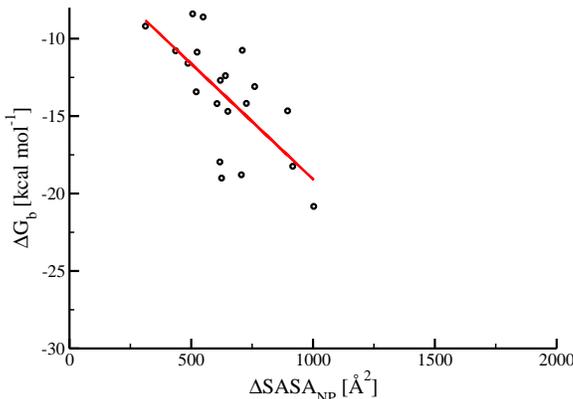
Every type of domain-domain interaction in iPfam is derived from at least one protein structure in which residues from the two domains are in close proximity to one another [21]. From that list of structures, we extracted a set of structures (PDB file names) representing every domain interaction occurring in the cSIN2. Given the thermodynamic differences governing intra- vs intermolecular interactions (see section 6.1 above), we removed from this set any structure in which the interaction domains are found on the *same* chain in the PDB file. We thus obtained a set of (co-crystal) structures containing all intermolecular instances of any domain-domain interaction occurring in the cSIN2.

For each structure containing an intermolecular domain-domain interaction – say between domains l and m – we calculated the change in solvent-accessible, non-polar surface area ( $\Delta\text{SASA}_{\text{NP}}$ ) using the software package POPS [22]. We proceeded by creating 3 separate PDB

files: file #1 contains only the atoms (ATOM records) that belong to the residues of the first domain **1**, file #2 contains only the atoms that belong to the residues of the second domain **m**, and file #3 contains the atoms from both **1** and **m**. In cases with more than one copy of a particular chain in the coordinates (as is often the case with structures derived from NMR), we used the coordinates from the first chain mentioned in the PDB file. Next, we calculated the non-polar solvent-accessible surface area for each file separately using POPs. This area is marked as “hydrophobic” in the POPs output. We then calculated  $\Delta\text{SASA}_{NP}$  as the difference between the sum of these areas for each domain separately and the area for the domains combined:

$$\Delta\text{SASA}_{NP}(1, m) = \text{SASA}_{NP}(1) + \text{SASA}_{NP}(m) - \text{SASA}_{NP}(1 + m). \quad (15)$$

We calculated  $\Delta\text{SASA}_{NP}$  using definition 15 for all the intermolecular domain-domain interactions in the cSIN2.



**Figure 11. Solvent-accessible non-polar surface areas and free energies of binding.**

The graph depicts the relationship between the change in solvent-accessible non-polar surface area,  $\Delta\text{SASA}_{NP}$  as defined in equation 15, and the free energy of binding. The free energies and the structures on which we based  $\Delta\text{SASA}_{NP}$  were taken from Table 1 in [23] and  $\Delta\text{SASA}_{NP}$  is calculated using POPs [22]. The red line represents a least-squares linear regression of the data, yielding equation 16 with  $R^2 = 0.47$ .

Several studies have noted that  $\Delta\text{SASA}_{NP}$  is related to the free energy of binding  $\Delta G_b$  [23, 24]. To map  $\Delta\text{SASA}_{NP}$  into  $\Delta G_b$ , we used a compilation of 20 individual structures that contain interactions whose  $\Delta G_b$  had been measured [23]. In Figure 11, we plot the  $\Delta\text{SASA}_{NP}$  of these structures, calculated using equation 15, against their measured  $\Delta G_b$ . We found that the relationship between the free energy of binding and the associated change in solvent-accessible non-polar surface area is approximately linear, with an  $R^2$  value of 0.47. We converted the  $\Delta\text{SASA}_{NP}$  values obtained for the interactions in the cSIN2 using the best-fit linear equation:

$$\Delta G_b(1, m) = -0.015 \cdot \Delta\text{SASA}_{NP}(1, m) - 4.17 \quad (16)$$

with  $\Delta\text{SASA}_{\text{NP}}(l, m)$  given in  $\text{\AA}^2$  and  $\Delta G_b(l, m)$  in  $\text{kcal mol}^{-1}$ .

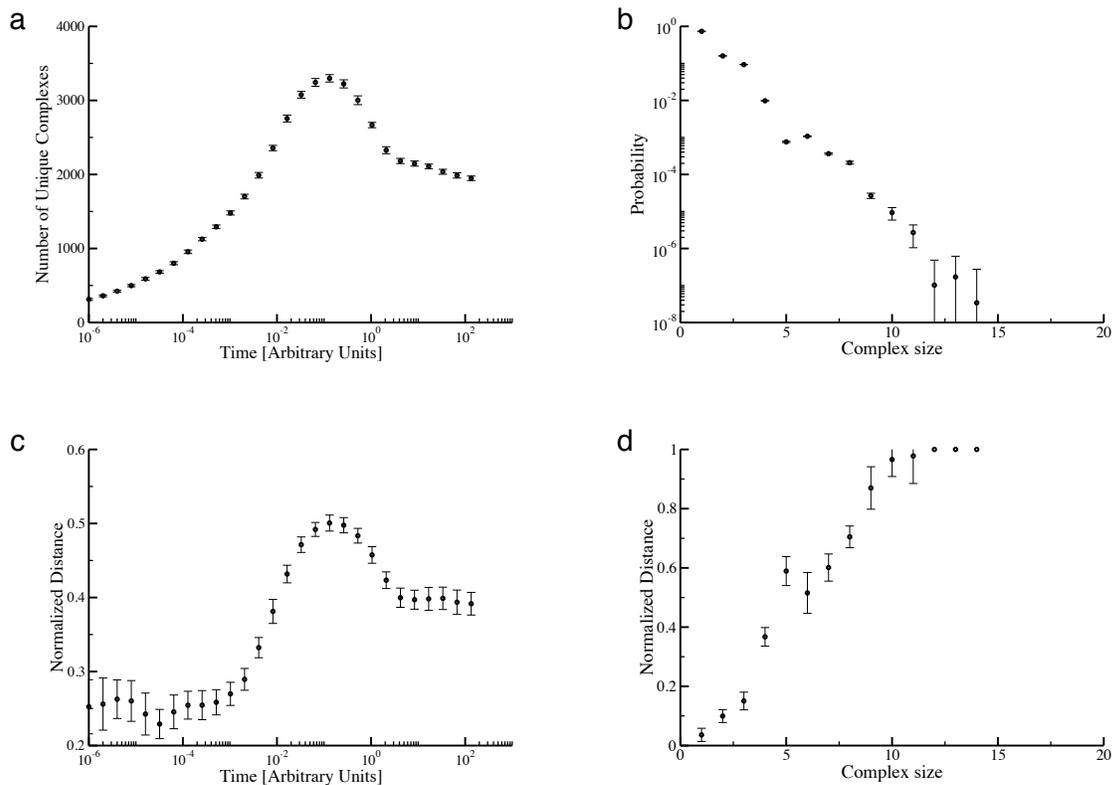
We used equation 16 to obtain binding free energies for the interactions in our cSIN2. In this process, we discarded any domain-domain interactions that buried less than  $100 \text{\AA}^2$  of non-polar surface area, as they are unlikely to contribute strongly to the overall protein-protein interaction as captured in that particular structure. The binding free energy of the interaction between domains  $l$  and  $m$  is then set to be the average  $\langle \Delta G_b(l, m) \rangle$  of the free energies resulting from all structures capturing the interaction between these domains and burying more than  $100 \text{\AA}^2$ . Finally, the  $K_D$  for each interaction is defined as  $K_D(l, m) = \exp\left(\frac{\langle \Delta G_b(l, m) \rangle}{RT}\right)$  with  $R$  the gas constant and  $T$  the absolute temperature. Here we consider interactions at room temperature, so  $RT$  is approximately  $0.6 \text{kcal mol}^{-1}$ .

### 8.3 Results for SR simulations

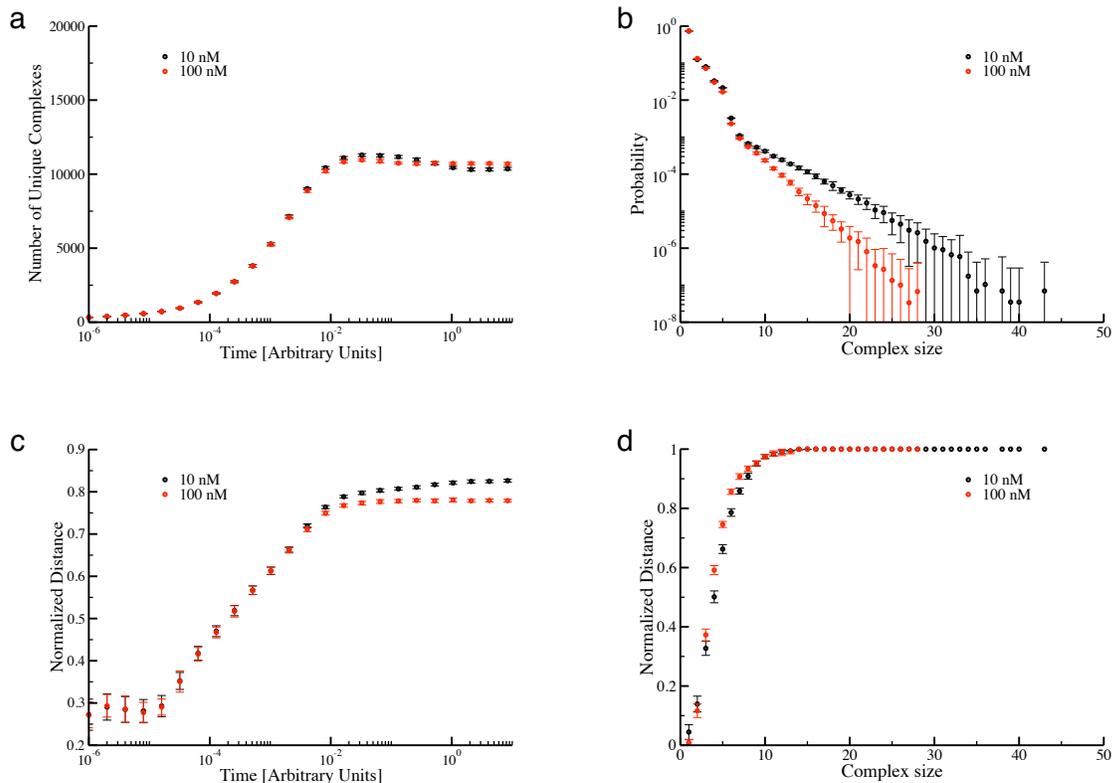
SR simulations behave quite differently from NR simulations, even for the same underlying parameter values, i.e. all  $K_D$  values set to 10 nM. The results for SR simulations at this affinity are shown in Figure 12, which is essentially analogous to Figure 5 (computed for the NR scenario) in the main text.

The SR constraint leads to considerably fewer unique molecular species at equilibrium (Figure 12a) when compared to the NR simulations ( $\sim 2000$  for the SR case vs.  $\sim 10000$  for the NR case, see Figure 5A in the main text). The system also takes much longer to reach equilibrium; the number of unique species initially climbs to a maximum value of  $\sim 3000$  but then decays slowly, reaching an apparent steady state after about 100 time units (compared to the  $\sim 8$  time units for the similar decay in the NR case, see Figure 5A in the main text). NR simulations that are run for the same length of time do not show any further signs of decay (data not shown). The SR scenario exhibits a more pronounced maximum because acyclic complexes that are formed initially eventually dissociate and are replaced by cyclic complexes, which are infinitely stable, resulting in a net decrease in the number of unique complexes on long time scales. Given the computational cost of running simulations that are this long (amounting to several weeks of CPU time), we cannot currently explore dynamics on a time scale longer than that shown in Figure 12a, which prevents us from determining whether the decay continues at a slow pace or has really ceased by  $\sim 100$  time units.

Figure 12b represents the distribution of complex sizes at the final time point in Figure 12a, analogous to Figure 6 in the main text for the NR scenario. Comparison of these figures reveals that SR simulations tend to sample much smaller complexes than NR simulations. SR simulations also tend to display less overall diversity; on long time scales the average distance between independent simulations is  $\sim 0.25$  (Figure 12c), as compared to  $\sim 0.7$  for the NR case (see Figure 5B in the main text). As can be seen from Figure 12d, however, SR simulations still exhibit significant diversity for large complexes. Indeed, the shape of this curve is essentially identical to that obtained in NR simulations (Figure 5C in the main text). The smaller overall distance between SR simulations is due to the fact that these simulations sample substantially fewer large complexes (Figure 12b), and does not indicate a significantly greater agreement between simulations in terms of which large complexes they contain. In fact, SR simulations tend



**Figure 12. Simulations with the SR constraint.** Results for simulations conducted using the SR constraints with all affinities set to  $K_D = 10$  nM. **(a)** This curve plots the average number of unique molecular species (averaged across 15 independent simulations) as a function of time. The error bars in this and all of the other panels in this figure represent  $\approx 95\%$  confidence intervals. The SR constraint clearly leads to many fewer unique complexes at steady state than the NR case. **(b)** The data points represent the average probability of finding a complex of a specified size across the entire population of 15 simulations at the final time point in panel **a**. The NR simulations have a much higher probability of sampling large complexes (i.e. complexes of size  $> 7$ ) than the SR simulations shown here. **(c)** This curve represents the distance between independent simulations, averaged across all unique comparisons between 5 simulations. As in panel **a**, the simulations exhibit a maximum distance ( $\sim 0.33$ ) before relaxing to a steady-state value of  $\sim 0.25$ . The distances observed between simulations are much smaller for the SR case than the NR case (compare with Figure 5B in the main text). **(d)** This curve represents the distance between simulations as a function of complex size, averaged over 5 simulations for the final time point in panel **a**. SR simulations show significant distances for large complexes. Taken together, this figure and panel **b** indicate that SR simulations exhibit smaller overall distances than NR simulations at steady state not because they exhibit greater similarity in their large complexes, but simply because they tend to sample much smaller complexes than NR simulations.



**Figure 13. Comparison of the 10 nM and 100 nM interaction affinity scenarios (NR constraint).** (a) Each curve represents the average number of unique complexes as a function of time (averaged over 15 independent simulations). The error bars in all panels in this figure represent  $\approx 95\%$  confidence intervals. Weakening the interaction affinity by one order does not strongly influence the steady state; the average number of unique complexes is essentially identical for the two cases. The 100 nM clearly lack the slight peak in the number of unique species observed in the 10 nM and actually exhibit slightly larger numbers of unique species at steady state. (b) Comparison of the distribution of complex sizes: the distributions represent the average probability of finding a complex of a particular size across the entire population of 15 simulations at the final time point in panel a. Networks with dissociation constant 10 nM have a higher probability of sampling large complexes (i.e. complexes of size  $> 7$ ) than those with dissociation constant 100 nM. (c) Comparison of the distance between independent simulations over time: each curve represents the average over all unique comparisons between 15 independent simulations using the distance measure defined in the main text. As in panel a, both interaction affinity scenarios produce strikingly similar curves; at steady state, the average normalized distance between  $K_D = 10$  nM simulations is  $\sim 0.83$  and for  $K_D = 100$  nM the average distance is  $\sim 0.78$ . (d) Comparison of the distance between independent simulations as a function of complex size: each curve represents the average over all unique comparisons between 15 independent simulations at the final time point in panel a. Again, the two affinity scenarios show very similar behavior; the major difference is that the 10 nM simulations sample much larger complexes than the 100 nM simulations.

to drift in the space of large complexes, much like NR simulations do (data not shown).

#### 8.4 Results for different affinity scenarios

In this section we present an overview of the results obtained for various interaction affinity scenarios in the context of both the SR and NR scenarios. The results shown and discussed in the main text are all obtained from NR simulations in which all affinities were set to  $K_D = 10$  nM. Figure 13 compares the 10 nM results to the 100 nM results for NR simulations; this figure is essentially analogous to Figure 5 in the main text.

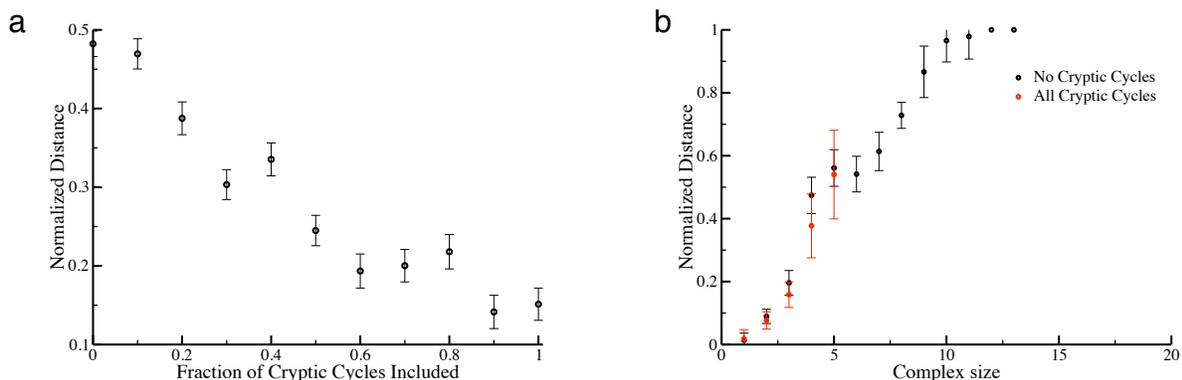
One finds that weaker interactions generally lead to smaller complexes (Figure 13b). This result is not surprising given that lower affinities give rise to fewer bonds and therefore smaller complexes in general. Since there are fewer unique options for smaller complexes (see Figure 4 in the main text), between-simulation distances are accordingly smaller (Figure 13c). As with the SR case at 10 nM, this is *not* due to the fact that simulations are *more* similar to each other with regard to the large complexes they contain, as demonstrated in Figure 13d. The smaller distances result from the fact that each simulation samples fewer large complexes, resulting in smaller overall distances even though the sets of large complexes in each simulation are essentially disjoint. We find the same general behavior for simulations at  $K_D = 1$   $\mu$ M. Such a low affinity sustains even fewer bonds, and thus even smaller complexes and inter-simulation distances (data not shown). Yet, even at this affinity, simulations disagree with regard to the large complexes they sample (data not shown).

#### 8.5 Results based on adding “cryptic” cycles

The thermodynamics of ring-like protein complexes (section 6.1) can result in situations in which a particular pair of sites might not bind one another strongly enough to be detected in a high-throughput interaction screen but could nonetheless contribute dramatically to the stability of certain complexes by forming a bond to complete a ring. Here we consider how such “cryptic” cycles (i.e. cycles that do not exist in the contact map but that may form very stable ring-like structures in a complex) influence the dynamics of our simulations.

By definition, the cSIN does not directly contain any information regarding cryptic cycles. In the absence of data on which proteins could form cryptic cycles, we add these cycles randomly to the graph and assess how increasing their number influences the behavior of the system. Given the computational cost of simulating the entire cSIN, we restricted this analysis just to the giant component of the cSIN (i.e. the cluster of proteins in the upper right-hand corner of Figure 2 in the main text). This allowed us to perform many simulations with varying amounts of cryptic cycles. For simplicity, we focused our analysis on cryptic cycles of length 3 (i.e. cycles containing only three proteins).

We first collected the set of protein triples that do not belong to a proper cycle in the giant component of the cSIN but could form a proper cycle if one edge were added to the graph. This set includes all proteins “A”, “B” and “C” such that (i) A, B and C are not part of a proper cycle in the cSIN, (ii) A can bind both B and C simultaneously, and (iii) both B and C have at least two



**Figure 14. Cryptic cycles.** Results based on simulations of the giant component of the cSIN in which cryptic cycles have been added to the graph. **(a)** The plot shows the average distance  $d$  for simulations of networks with a varying fraction of cryptic cycles added to the giant component of the cSIN, as described in the text. The “0” point corresponds to a graph in which no cryptic cycles are added (i.e. the giant component of the cSIN), while the point at a fraction of “1” represents a case in which all possible 113 cryptic cycles were added to the graph. The simulations were performed using the SR constraint and a constant  $K_D$  of 10 nM for all original reactions in the network (i.e. *not* the reactions added to generate cryptic cycles). Each point represents the average distance over all unique comparisons between 10 independent simulations at a particular time  $t$  after the simulations reach steady state (taken to be  $\sim 30$  time units, see Figure 12c). The error bars represent  $\approx 95\%$  confidence intervals. The average steady-state distance decreases quite dramatically as the number of cryptic cycles increases. **(b)** Here we plot the distance between complexes as a function of complex size for simulations with “No Cryptic Cycles” (the point corresponding to a value of 0 on the abscissa of panel a) and for simulations including “All Cryptic Cycles” (corresponding to the point at 1 in panel a). As one can see, adding cryptic cycles essentially reduces the maximum complex size observed in the simulation (in this case from 13 with no cycles to 5 when all possible cycles are included). The reduction in overall distance observed at higher cryptic cycle densities results from these cycles effectively preventing the sampling of large complexes.

sites. In the giant component of the cSIN we identified 113 distinct sets of proteins meeting the criterion of a “potential” cryptic cycle. We then created new versions of the giant component of the cSIN in which a fraction (20%, 30%, etc.) of these potential cycles were converted into proper cycles by adding a corresponding binding rule to the system. Given the underlying hypothesis of the cryptic cycle (i.e. that binding between B and C is not strong enough to occur appreciably on its own), we added these new rules with very low on-rates ( $\beta_+ = 10^{-8}$ ). In SR simulations cryptic cycles will still form infinitely stable ring-like structures despite the low on-rate (low affinity) of the B-C bond (see section 6.2).

We constructed a set of graphs with a fraction of cryptic cycles ranging from 10% to 100%. In each graph, the set of cryptic cycles added to the system was chosen uniformly at random from the set of 113 possibilities. If more than one pair of sites could be used to generate a proper cycle from a given potential cycle, we chose one such pair at random with equal probability from the set of possible pairs. We ran 10 independent simulations for each graph and calculated the average steady-state distance between simulations (using distance “ $d$ ” from equation 11) as a function of the fraction of cryptic cycles included in the graph. The results are shown in Figure 14.

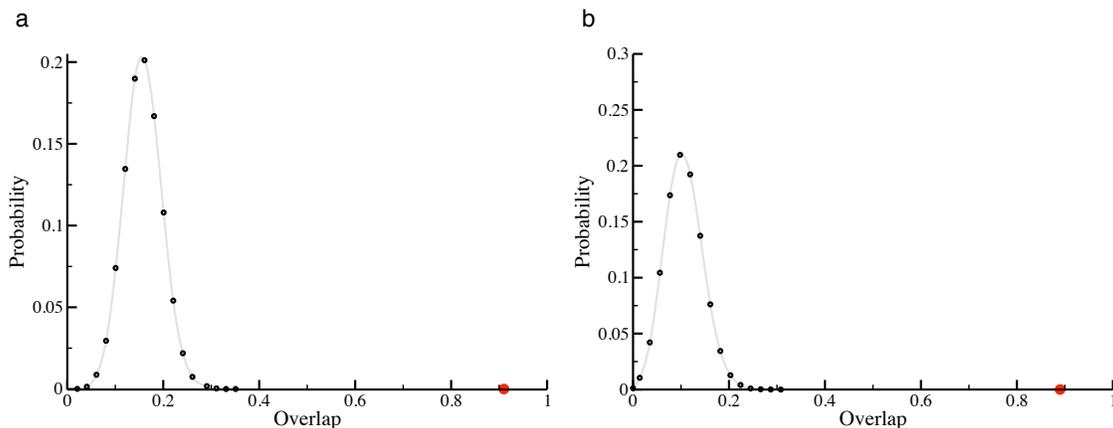
Adding cryptic cycles to the graph clearly decreases the overall distance between simulations at steady state. Including all of the possible cycles (corresponding to a fraction of 1 in Figure 14a) results in a distance of only 0.15 compared to a distance of 0.5 when no cryptic cycles are added. The addition of cryptic cycles results in simulations that tend to sample much smaller complexes: when all cryptic cycles are added, the largest complex formed by the system has only 5 members, compared to 13 when no cryptic cycles are present. The net effect of including cryptic cycles is much like the effect of including the SR (vs. NR) constraint; the overall distance between cells decreases not because cryptic cycles increase the similarity between large complexes but rather because cryptic cycles dramatically reduce the size of complexes sampled by the simulation.

## 9 Comparison with Affinity Purification / Mass Spectrometry data

In this section we aim at validating the structure of the cSIN (beyond its original curation by Kim *et al* [19]) and the soundness of our dynamical process (including global constraints) by comparing simulation outcomes with available data germane to complex formation in yeast cells. At present, there is no direct experimental method that can determine all of the individual complexes present in a single cell at a given time, but we can compare simulation outcomes to Affinity Purification-Mass Spectrometry (AP-MS) data [25, 26]. This type of data samples complexes from a large number of cells and provides only limited insight into the detailed structure of the complexes themselves. Although one can attempt to reconstruct sets of complexes from this data (e.g. [26, 27]), many different sets of putative complexes can be equally consistent with a given set of empirical results.

AP-MS experiments ultimately produce putative co-complex relationships [27]. Our simulations produce a set of structured complexes that we can convert into a co-complex relationship of the AP-MS kind by defining two proteins A and B to be related if they are found together in at least one complex in a population of simulations at a particular time  $t$ . We create our set of co-complex relationships by examining the 15 independent simulations for a given set of conditions (either the SR or NR scenarios at  $K_D = 10$  nM) at a particular time point after steady-state in the number of unique complexes per cell has been reached ( $\sim 8$  time units for the NR scenario and  $\sim 100$  time units for the SR scenario).

We then compare the set of simulated co-complex interactions to the set observed in a recent collation of AP-MS data employed by Yu *et al.* [25]. This data was downloaded directly from the authors’ website at [http://interactome.dfci.harvard.edu/S\\_cerevisiae](http://interactome.dfci.harvard.edu/S_cerevisiae). We calculate a parameter  $f$  that represents the fraction of interactions observed in the data that are also observed in our simulation. Naturally,  $f$  is restricted to the set of interactions between proteins that are included in the cSIN, which yields a total of 48 co-complex interactions in the data. Of these 48



**Figure 15. Overlap with AP-MS data.** Distribution of the overlap between randomized SR co-complex interactions and the AP-MS data from [25]. Panel (b): Simulations with the SR scenario. In both of these plots we compare the outcomes of our simulations with AP-MS data. The abscissa indicates the overlap between simulation data and co-complex interactions observed by AP-MS. The red dot marks the overlap between our simulations and experimental data. The black data points represent a histogram of the frequency with which a particular overlap is observed in a set of 100,000 randomizations of our simulation data (“randomized overlap”), as described in the text. The grey curve represents a smoothed version of this distribution. The separation between the randomized overlap and the actual overlap indicates that our simulation outcomes are extremely significant. Note: This plot is based on a limited set of interactions (only 48 interactions involving cSIN proteins are actually covered in the AP-MS data). Panel (a): Simulations with the NR scenario.

interactions, however, 11 occur between pairs of proteins that belong to different components in the cSIN (i.e. are found in different clusters in Figure 2 of the main text). Since there is no path in the graph linking the two proteins in these 11 pairs, our simulation cannot produce co-complex interactions for any of them.

As discussed in the main text, the cSIN is derived from the giant component of the SIN, indicating that the underlying PPI data contains a path connecting each pair in this set of 11. In most cases, these paths are missing from the cSIN due to the fact that proteins on this path had either “ambiguous” localizations or were simply not visualized in the experimental localization data [28]. In a few cases, a path exists between the two proteins in question in a different compartment (such as the nucleus). Since AP-MS experiments assay all of the compartments in the cell, they recover co-complex interactions that we cannot observe in simulations of the cytoplasm.

If we restrict our comparison only to those interactions that do not suffer from the localization issues discussed above, we obtain an  $f$  of 91% for the NR simulations and 89% for the SR simulations. Considering all of the 48 co-complex interaction possibilities reduces the overlap to 71% and 69%, respectively.

To understand whether the degree of overlap we observe could be explained by random

association of proteins in a co-complex graph, we performed a randomization of each set of simulated co-complex interactions by randomly swapping interaction pairs (essentially randomly rewiring the co-complex graph). For both the SR and the NR scenario, we performed  $10^5$  such randomizations. The comparison between the distribution of  $f$ -values in these randomized sets and the actual simulations is shown in Figure 15 for the NR and SR scenarios. In both scenarios we see that the p-value of observing the values of  $f$  that we see is  $\ll 10^{-5}$ ; this is true regardless of whether we consider the simulation overlaps to be  $\sim 70\%$  rather than  $\sim 90\%$ . We have also compared our results to the set of “core” complexes observed by Gavin et al. [26], with essentially identical results to those shown above (data not shown).

We found that the level of overlap for NR and SR simulations did not vary with interaction strength, nor did either case produce a different overlap with experiment when affinities were determined using the concentration-based scenario (derived from equation 4 in the main text). The cSIN2, however, results in only a 54% overlap with AP-MS results. This is because a number of edges in the cSIN are missing from the cSIN2, as suitable structures to define the affinities between certain types of domain-domain interactions could not be found. These edges evidently carry significant levels of information from the standpoint of complex formation (given the much lower overlap between cSIN2 simulations and the experimental results). For this reason we focus on the 10 nM results in the main text, although the overall behavior of the cSIN2 is very similar to the constant-affinity case of the original cSIN (see Figure 9 in the main text).

It is important to note two facts about this comparison. First, the AP-MS data from Yu et al. contain only a small number of distinct co-complex interactions when restricted to cSIN proteins. Indeed, about 50% of the proteins included in our simulations are not mentioned in this data set at all. Thus, the values of  $f$  calculated above, while quite statistically significant, are based on a comparatively small number of data points. Second, our simulations predict hundreds of co-complex associations that are not observed, yet *could* be detectable in principle (because they occur between two proteins that are both mentioned in the AP-MS dataset; for which, therefore, some number of co-complexes have been identified). Many such interactions are actually purely binary interactions observed in highly curated datasets (including the SIN) yet are not found in the AP-MS data [25]. Either these interactions do not exist in cells (and thus represent errors in our approach) or they do exist but have not been observed (and thus represent errors or omissions in the AP-MS data).

## References

1. Blinov ML, Faeder JR, Hlavacek WS (2004) BioNetGen: Software for rule-based modeling of signal transduction based on the interactions of molecular domains. *Bioinformatics* 20: 3289-3292.
2. Danos V, Laneve C (2004) Formal molecular biology. *Theoretical Computer Science* 325: 69-110.
3. Lok L, Brent R (2005) Automatic generation of cellular reaction networks with molecuizer 1.0. *Nature Biotechnology* 23: 131-136.

4. Danos V, Feret J, Fontana W, Harmer R, Krivine J (2007) Rule-based modelling of cellular signalling. In: Proceedings of the 18th Int. Conf. on Concurrency Theory. Lisboa, Portugal: Springer, volume 4703 of *Lecture Notes in Computer Science*, pp. 17–41.
5. Mallavarapu A, Thomson M, Ullian B, Gunawardena J (2008) Programming with models: modularity and abstraction provide powerful capabilities for systems biology. *J Roy Soc Interface* doi:10.1098/rsif.2008.0205.
6. Hlavacek WS, Faeder JR, Blinov ML, Posner RG, Hucka M, et al. (2006) Rules for modeling signal-transduction systems. *Science STKE* 344: re6.
7. Feret J, Danos V, Krivine J, Harmer R, Fontana W (2009) Internal coarse-graining of molecular systems. *Proc Natl Acad Sci USA* 106: 6453–6458.
8. Harmer R, Danos V, Feret J, Krivine J, Fontana W (2010) Intrinsic information carriers in combinatorial dynamical systems. *Chaos* 20: 037108.
9. Doob JL (1945) Markoff chains denumerable case. *Trans Amer Math Soc* 58: 455–473.
10. Gillespie DT (1976) A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of Computational Physics* 22: 403–434.
11. Danos V, Feret J, Fontana W, Krivine J (2007) Scalable simulation of cellular signalling networks. In: Proceedings APLAS 2007. Springer, volume 4807 of *Lecture Notes in Computer Science*, pp. 139–157.
12. Yang J, Monine MI, Faeder JR, Hlavacek WS (2008) Kinetic monte carlo method for rule-based modeling of biochemical networks. *Phys Rev E* 78: 031910.
13. Danos V, Feret J, Fontana W, Krivine J (2008) Abstract interpretation of cellular signalling networks. In: Verification, Model Checking, and Abstract Interpretation. Springer, volume 4905 of *Lecture Notes in Computer Science*, pp. 83–97.
14. Wolfram Research, Inc (2008) Mathematica Edition: Version 7.0. Champaign, Illinois: Wolfram Research, Inc.
15. Saiz L, Vilar JM (2006) Stochastic dynamics of macromolecular-assembly networks. *Mol Syst Biol* 2: 2006 0024.
16. Minh DD, Bui JM, Chang CE, Jain T, Swanson JM, et al. (2005) The entropic cost of protein-protein association: a case study on acetylcholinesterase binding to fasciculin-2. *Biophys J* 89: L25–7.
17. Bray D, Lay S (1997) Computer-based analysis of the binding steps in protein complex formation. *Proc Natl Acad Sci U S A* 94: 13493–8.
18. Jorgensen P, Nishikawa JL, Breitzkreutz BJ, Tyers M (2002) Systematic identification of pathways that couple cell growth and division in yeast. *Science* 297: 395–400.
19. Kim PM, Lu LJ, Xia Y, Gerstein MB (2006) Relating three-dimensional structures to protein networks provides evolutionary insights. *Science* 314: 1938–41.

20. Finn RD, Tate J, Mistry J, Coghill PC, Sammut SJ, et al. (2008) The pfam protein families database. *Nucleic Acids Res* 36: D281–8.
21. Finn R, Marshall M, Bateman A (2004) ipfam: visualization of protein-protein interactions in pdb at domain and amino acid resolutions. *Bioinformatics* 21: 410-412.
22. Fraternali F, Cavallo L (2002) Parameter optimized surfaces (pops): analysis of key interactions and conformational changes in the ribosome. *Nucleic Acids Res* 30: 2950-2960.
23. Bougouffa S, Warwicker J (2008) Volume-based solvation models out-perform area-based models in combined studies of wild-type and mutated protein-protein interfaces. *BMC Bioinformatics* 9: 448.
24. Horton N, Lewis M (1992) Calculation of the free energy of association for protein complexes. *Protein Sci* 1: 169-81.
25. Yu H, Braun P, Yildirim MA, Lemmens I, Venkatesan K, et al. (2008) High-quality binary protein interaction map of the yeast interactome network. *Science* 322: 104–110.
26. Gavin AC, Aloy P, Grandi P, Krause R, Boesche M, et al. (2006) Proteome survey reveals modularity of the yeast cell machinery. *Nature* 440: 631–636.
27. Scholtens D, Gentleman R (2004) Making sense of high-throughput protein-protein interaction data. *Stat Appl Genet Mol Biol* 3: Article39.
28. Huh WK, Falvo JV, Gerke LC, Carroll AS, Howson RW, et al. (2003) Global analysis of protein localization in budding yeast. *Nature* 425: 686–91.