Supplementary Notes

Network Deconvolution - A General Method to Distinguish Direct Dependencies over Networks

Soheil Feizi^{1,2,3,4}, Daniel Marbach^{1,2}, Muriel Médard³ and Manolis Kellis^{1,2,4}

¹ Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, USA.

² Broad Institute of MIT and Harvard, Cambridge, MA, USA.

³ Research Laboratory of Electronics (RLE) at MIT, Cambridge, Massachusetts, USA.

⁴ Correspondent author (sfeizi@mit.edu,manoli@mit.edu).

Supplementary Sections

1	Analy	yzing General Properties of Network Deconvolution	3
	1.1	Network deconvolution algorithm	3
	1.2	Optimality analysis of network deconvolution	3
	1.3	Modeling assumptions and intuitions	4
	1.4	Decomposable and non-decomposable matrices	4
		1.4.1 Asymmetric decomposable matrices	4
		1.4.2 General non-decomposable matrices	5
	1.5	Robustness against noise	6
	1.6	Scaling effects	6
	1.7	Computational complexity analysis of network deconvolution	7
		1.7.1 Using sparsity in eigen decomposition	7
		1.7.2 Parallelization of network deconvolution	8
2	Inferm	ring Gene Regulatory Networks by Network Deconvolution	8
	2.1	DREAM framework	8
	2.2	Mutual Information	9
	2.3	Evaluation functions	9
	2.4	Network motif analysis	9
3	Inferr	ring Protein Structural Constraints by Network Deconvolution	10
4	Inferr	ring Weak and Strong Collaboration Ties by Network Deconvolution	10

Figures

S1	Effects of applying network deconvolution on network eigenvalues	12
S2	Convergence of network deconvolution for non-decomposable matrices	13
S3	Robustness of network deconvolution against noise	14
S4	Linear scaling effects on network deconvolution performance	15

S5	Scalability of network deconvolution to large networks using eigen sparsity and parallelization	16
S6	DREAM framework	17
S7	Overall total and non-redundant contact discovery rates	18
$\mathbf{S8}$	Discovery rate of interacting contacts for all tested proteins	19
$\mathbf{S9}$	Discovery rate of non-redundant contacts for all tested proteins	20
S10	Total and non-redundant discovery rate for different thresholds of contact proximity $\ldots \ldots \ldots$	21

Tables

S1	General properties of considered networks	22
S2	Tested protein names and PDB IDs	23

1 Analyzing General Properties of Network Deconvolution

In this section, we describe the proposed *network deconvolution* (ND) algorithm and analyze its general properties. First, we show how network deconvolution finds a globally optimal direct dependency matrix by eliminating indirect dependencies from the observed network. We present an extension of network deconvolution algorithm for non-decomposable matrices based on an iterative method. Then, we analyze the robustness of network deconvolution in the presence of noise and linear scaling. Finally, we analyze computational complexity of network deconvolution and propose two extensions to it (using eigen sparsity and parallelization) to be able to efficiently use network deconvolution over very large networks.

1.1 Network deconvolution algorithm

Network deconvolution is a systematic approach of computing direct dependencies in a network by use of local edge weights. Recall that, G_{obs} represents an observed dependency matrix, a properly scaled similarity matrix among variables. The linear scaling depends on the largest absolute eigenvalue of the un-scaled similarity matrix and is discussed in more details in Section 1.6. Components of G_{obs} can be derived by use of different pairwise similarity metrics such as correlation or mutual information. In particular, $g_{i,j}^{obs}$, the (i, j)-th component of G_{obs} , represents the similarity value between the *observed* patterns of variables *i* and *j* in the network.

A perennial challenge to inferring networks is that, both direct and indirect dependencies arise together. A direct information flow modeled by an edge in G_{dir} can give rise to two or higher level indirect information flows captured in G_{indir} :

$$G_{indir} = G_{dir}^2 + G_{dir}^3 + \dots$$
 (1)

The power associated with each term in G_{indir} corresponds to the level of indirection contributed by that term. We assume that the observed dependency matrix, G_{obs} , comprises both direct and indirect dependency effects (i.e., $G_{obs} = G_{dir} + G_{indir}$). The case of having external noise over observed dependencies is considered in Section 1.5. Further, we assume that G_{dir} is an $n \times n$ decomposable matrix. The case of non-decomposable matrices is considered in Section 1.4.

In Section 1.2, we show that the following network deconvolution algorithm finds an optimal solution for direct dependency weights by using the observed dependencies:

Algorithm 1 Network deconvolution has three steps:

- Linear Scaling Step: The observed dependency matrix is scaled linearly so that all eigenvalues of the direct dependency matrix are between -1 and 1.
- Decomposition Step: The observed dependency matrix G_{obs} is decomposed to its eigenvalues and eigenvectors such that $G_{obs} = U \Sigma_{obs} U^{-1}$.
- Deconvolution step: A diagonal eigenvalue matrix Σ_{dir} is formed whose *i*-th component is $\lambda_i^{dir} = \frac{\lambda_i^{obs}}{\lambda_i^{obs}+1}$. Then, the output direct dependency matrix is $G_{dir} = U\Sigma_{dir}U^{-1}$.

1.2 Optimality analysis of network deconvolution

In this section, we show how the network deconvolution algorithm proposed in Algorithm 1 finds an optimal solution for direct dependency weights by using the observed ones.

Suppose U and Σ_{dir} represent eigenvectors and a diagonal matrix of eigenvalues of G_{dir} , where λ_i^{dir} is the *i*-th diagonal component of the matrix Σ_{dir} . Also, suppose the largest absolute eigenvalue of G_{dir} is strictly smaller than one. This assumption holds by using a linear scaling function over the un-scaled observed dependency network and is discussed in Section 1.6. By using the eigen decomposition principle⁷, we have $G_{dir} = U\Sigma_{dir}U^{-1}$. Therefore,

$$\begin{aligned} G_{dir} + G_{indir} &\stackrel{(a)}{=} & G_{dir} + G_{dir}^2 + \dots \qquad (2) \\ \stackrel{(b)}{=} & (U\Sigma_{dir}U^{-1}) + (U\Sigma_{dir}^2U^{-1}) + \dots \\ &= & U(\Sigma_{dir} + \Sigma_{dir}^2 + \dots)U^{-1} \\ &= & U\left(\begin{array}{cc} \sum_{i\geq 1}(\lambda_1^{dir})^i & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sum_{i\geq 1}(\lambda_n^{dir})^i \end{array}\right)U^{-1} \\ \stackrel{(c)}{=} & U\left(\begin{array}{cc} \frac{\lambda_1^{dir}}{1-\lambda_1^{dir}} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{\lambda_n^{dir}}{1-\lambda_d^{dir}} \end{array}\right)U^{-1}. \end{aligned}$$

Equality (a) follows from the definition of diffusion model of equation (1). Equality (b) follows from the eigen decomposition of matrix G_{dir} . Equality (c) uses geometric series to compute the infinite summation in a closed-form since $|\lambda_i^{dir}| < 1$ for all *i*.

By using the eigen decomposition of the observed network, G_{obs} , we have $G_{obs} = U\Sigma_{obs}U^{-1}$, where

$$\Sigma_{obs} = \begin{pmatrix} \lambda_1^{obs} & \cdots & 0\\ \vdots & \ddots & \vdots\\ 0 & \cdots & \lambda_n^{obs} \end{pmatrix}.$$
 (3)

Therefore, from equations (2) and (3), if

$$\frac{\lambda_i^{dir}}{1 - \lambda_i^{dir}} = \lambda_i^{obs} \quad \forall 1 \le i \le n,$$
(4)

the error term $G_{obs} - (G_{dir} + G_{indir}) = 0$. Rewriting equation (4) leads to a non-linear filter over eigenvalues in the network deconvolution algorithm:

$$\lambda_i^{dir} = \frac{\lambda_i^{obs}}{1 + \lambda_i^{obs}} \quad \forall 1 \le i \le n.$$
(5)

This nonlinear filter over eigenvalues allows network deconvolution to eliminate indirect dependency effects modeled as in equation (1) and compute a globally optimal solution for direct dependencies over the network with zero error.

1.3 Modeling assumptions and intuitions

In this section, we provide more intuitions on network deconvolution and discuss about its modeling assumptions and limitations.

As discussed in Sections 1.1 and 1.2, network deconvolution can be viewed as a nonlinear filter that is applied on eigenvalues of the observed dependency matrix. Figure S1-a shows an example of how eigenvalues of the observed matrix change by application of network deconvolution. In general, ND decreases magnitudes of large positive eigenvalues of the observed matrix since transitivity effects modeled as in equation (1) inflate these positive eigenvalues. Network deconvolution reverses this effect by using the described nonlinear filter.

There are some modeling assumptions and limitations of network deconvolution as follows:

- Network deconvolution assumes that networks are *linear time-invariant flow-preserving* operators. Under this condition, indirect flows over different paths are modeled by using power series of the direct matrix, while observed dependencies are modeled as the sum of direct and indirect effects.
- We assume that the maximum absolute eigenvalue of the direct dependency matrix is strictly smaller

than one. This assumption can be satisfied by linearly scaling the observed dependency matrix. This is explained in details in Section 1.6.

When flow assumptions of network deconvolution hold, network deconvolution removes all indirect flow effects and infers all direct interactions and weights exactly (**Figure 1d**). When these assumptions do not hold by inclusion of non-linear probabilistic effects through simulations, we show that ND performs well and infers most of direct edges (**Figure 1e**). In our simulations, we use scale-free networks with 100 nodes and add transitive edges using a nonlinear probabilistic function considering two-hop and three-hop indirect paths over the network. Finally, note that if the observed network is a covariance matrix of jointly Gaussian variables, ND infers direct interactions using global partial correlations.

1.4 Decomposable and non-decomposable matrices

The network deconvolution framework described in Algorithm 1 is based on eigen decomposition principles and therefore only can be applied on so called *diagonalizable* matrices⁷: a matrix G is diagonalizable if it can be decomposed to its eigenvectors U and a diagonal matrix of eigenvalues Σ such that $G = U\Sigma U^{-1}$. Many symmetric matrices and some asymmetric matrices are diagonalizable⁷.

In this section, first we present an example of an asymmetric matrix which is diagonalizable and has real eigenvalues. Then, we propose an extension of network deconvolution for general non-diagonalizable matrices based on an iterative conjugate gradient method.

1.4.1 Asymmetric decomposable matrices

Here, we present an example of an asymmetric matrix which is decomposable to its eigenvectors and real eigenvalues. The structure of this asymmetric matrix is used in the application of network deconvolution in inferring gene regulatory networks where the observed dependency matrix is asymmetric.

Example 2 Suppose G is a $n \times n$ matrix with the following structure:

$$G = \begin{bmatrix} G_1 & G_2 \\ \hline 0 & 0 \end{bmatrix},\tag{6}$$

where G_1 is a $m \times m$ symmetric matrix and G_2 is a $m \times (n-m)$ matrix which can be asymmetric. We show that, eigenvalues of G are real.

Suppose λ is an eigenvalue of G which corresponds to the eigenvector v. By definition, $Gv = \lambda v^7$. Say

$$v = [v_1 \mid v_2]',$$

where v_1 contains the first m elements of v. Therefore,

$$Gv = \lambda v \Rightarrow \left[\frac{G_1v_1 + G_2v_2}{0}\right] = \lambda \left[\frac{v_1}{v_2}\right].$$
 (7)

From equation (7), we have $v_2 = 0$ and $G_1v_1 = \lambda v_1$. In other words, λ is also an eigenvalue for the matrix G_1 . Since the matrix G_1 is symmetric and diagonalizable, λ is real and G is diagonalizable.

Example 2 illustrates a matrix structure which is common in gene regulatory network inference problems. In that application, the first m rows correspond to transcription factor genes and hence they can have outgoing and incoming edges. The last n - m rows correspond to genes that are not transcription factors and therefore only have incoming edges. This example illustrates that the proposed network deconvolution algorithm based on eigen decomposition principles can be applied in gene regulatory network inference problems. We discuss more details of this application of network deconvolution in Section 2.

1.4.2 General non-decomposable matrices

Most symmetric and some asymmetric matrices are decomposable and hence amenable to the network deconvolution algorithm based on eigen decomposition principles. However, if the input matrix is not decomposable, an extension of network deconvolution based on an iterative gradient descent method can be used to eliminate indirect information effects from observed dependencies. We show that, for a general asymmetric matrix, this variation of ND finds a direct dependency matrix in a network. In the case of convexity, the proposed algorithm converges to a globally optimal solution. We also illustrate the convergence of this iterative algorithm over a simulated asymmetric network.

For simplicity, in this part, we only consider the second degree term of the diffusion model (i.e., $\hat{G}_{indir} = G^2$). If the largest absolute eigenvalue of G_{dir} is small enough, higher order diffusion terms are decaying exponentially and ignoring them can be justified. Also, suppose the observed dependency matrix G_{obs} is scaled so that $g_{i,j}^{obs} \leq 1/2$. This assumption is required to have a convex optimization setup and can be obtained by linearly scaling the observed dependency matrix.

Suppose $\Gamma(G_{dir})$ represents the energy of the error term:

$$\Gamma(G_{dir}) \triangleq \|G_{obs} - (G_{dir} + \hat{G}_{indir})\|_{F}^{2}$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{n} (g_{i,j}^{obs} - (g_{i,j}^{dir} + \hat{g}_{i,j}^{indir}(G_{dir})))^{2}, (8)$$

where $\|.\|_F$ represents the Frobenius norm⁷. To compute G_{dir} , the energy of the error term is minimized. This can be formulated as the following optimization problem:

$$\min_{G_{dir}} \Gamma(G_{dir}). \tag{9}$$

Suppose $\nabla \Gamma(G_{dir})$ is the gradient of the matrix $\Gamma(G_{dir})$ with respect to G_{dir} . The (i, j)-th component of $\nabla \Gamma(G_{dir})$ is:

$$[\nabla\Gamma]_{i,j} = \frac{\partial\Gamma}{\partial g_{i,j}^{dir}}.$$
(10)

A closed form expression for $\nabla \Gamma(G_{dir})$ can be derived by derivative calculations.

Algorithm 3 The following gradient descent algorithm finds a globally optimal solution for the optimization problem of equation (9):

- Step 0: $G^1_{dir} = G_{obs}$.
- Step i: $G_{dir}^{i+1} = G_{dir}^i \beta_i \nabla \Gamma(G_{dir}^i)$.

Step sizes β_i is a decreasing sequence and can be chosen in various ways². Since the optimization problem of equation (9) is convex (the condition $g_{i,j}^{obs} \leq 1/2$ is required for this), this algorithm converges to a globally optimal solution². In non-convex cases, the proposed gradient descent algorithm converges to a locally optimal solution.

We illustrate the convergence of Algorithm 3 over a randomly generated asymmetric (directed) network with 50 nodes and 490 edges. Figure S2 illustrates how the energy of error goes to zero as the number of iterations increases.

Using the iterative conjugate gradient method can be very expensive especially for very large networks. For example, in the worst case time complexity, at each iteration, we need to compute n^2 gradient terms of equation (10), where each term can be computed in $\mathcal{O}(n)$ worst case time complexity. Therefore, for a fixed number of interactions, the worst case time complexity of the iterative algorithm can be as high as $\mathcal{O}(n^3)$. However, we have developed two extensions of network deconvolution that make it scalable to very large networks (Section 1.7): the first exploits the sparsity of eigenvalues of low rank networks, and the second parallelizes network deconvolution over potentiallyoverlapping subgraphs of the network.

1.5 Robustness against noise

In this section, we show that network deconvolution is provably robust against noise. We also illustrate its robustness in the application of ND in gene regulatory network inference.

Suppose the observed dependency matrix is itself noisy:

$$\hat{G}_{obs} = G_{obs} + N = G_{dir} + G_{indir} + N,$$

where N represents the noise matrix. In this case, we may not be able to recover G_{dir} exactly from \hat{G}_{obs} but instead provide an estimate for it, referred to as \hat{G}_{dir} , and show that, this estimate is close to G_{dir} as long as the noise is moderate.

We assume that noise is moderate. We quantify the noise power by its Euclidean norm, i.e., $\sqrt{(\sum_{i,j} n_{i,j}^2)}$ where $n_{i,j}$ is the (i, j)-th component of the noise matrix. This norm can be calculated⁷ by the largest absolute eigenvalue of the noise matrix N. Here, we assume that the largest absolute eigenvalue of N, referred to as γ , is much smaller than 1 ($\gamma << 1$), leading to moderate noise. Also, we assume that the largest absolute eigenvalue of \hat{G}_{obs} is equal to $\delta <$ 1, leading to convergence of the Taylor series. This can be obtained by using a linear mapping function over unscaled observed dependencies. Thus, when the power of the noise is smaller than the power of the observed network ($\gamma << 1$ and $\delta < 1$), we have sufficient information for robust network deconvolution.

By use of network deconvolution over the noisy observed network \hat{G}_{obs} , an estimate of direct dependency matrix \hat{G}_{dir} is obtained. We show that, this estimate is close to G_{dir} by bounding the Euclidean norm of the error term (i.e., $\|G_{dir} - \hat{G}_{dir}\|_{l_2}$) as a function of δ and γ , the largest absolute eigenvalues of the observed dependency network \hat{G}_{obs} and the noise matrix N, respectively:

$$\begin{split} G_{dir} - \hat{G}_{dir} \|_{l_{2}} &\stackrel{\text{(d)}}{=} \| \frac{G_{obs}}{I + G_{obs}} - \frac{\hat{G}_{obs}}{I + \hat{G}_{obs}} \|_{l_{2}} \\ &\stackrel{\text{(e)}}{=} \| (G_{obs} - (G_{obs})^{2} + \ldots) \\ &- (\hat{G}_{obs} - (\hat{G}_{obs})^{2} + \ldots) \|_{l_{2}} \\ &= \| (G_{obs} - \hat{G}_{obs}) \\ &- (G_{obs}^{2} - (\hat{G}_{obs})^{2}) + \ldots \|_{l_{2}} \\ &\stackrel{\text{(f)}}{\leq} \| G_{obs} - \hat{G}_{obs} \|_{l_{2}} \\ &+ \| G_{obs}^{2} - \hat{G}_{obs}^{2} \|_{l_{2}} + \ldots \\ &\stackrel{\text{(g)}}{\leq} \gamma + \gamma^{2} + 2\delta\gamma + \ldots \\ &\stackrel{\text{(h)}}{\leq} \gamma + \mathcal{O}(\delta^{2} + \gamma^{2} + \delta\gamma). \end{split}$$
(11)

Equality (d) follows from the application of network deconvolution over matrices G_{obs} and \hat{G}_{obs} as in equation (5). Equality (e) follows from the Taylor series expansion of the function $\frac{x}{1+x}$ for matrices. Inequalities (f),(g) and (h) follow from basic algebraic inequalities⁷.

Equation (11) shows that, for small enough δ , network deconvolution is robust against noise. Hence, in different applications, we linearly scale the observed network to have $\delta < 1$. The effect of this linear mapping is discussed in more details in Section 1.6.

Finally, we illustrate the robustness of network deconvolution in the gene regulatory network inference problem. An inference method is robust against noise when its performance degrades *continuously* by increasing the noise power. Here, for each experiment, we add an artificial Gaussian noise with a variance of a fraction of the experiment variance and use these noisy datasets to infer networks. Figure S3 shows the performance of network deconvolution algorithm in different noise levels in *E. coli* and *in silico* networks. Here, to observe noise effects completely, unlike the DREAM framework scoring (see Section 2.3), we do not threshold inferred networks. As it can be seen in this figure, by increasing the noise power, the overall score decreases smoothly which indicates the robustness of network deconvolution against noise.

1.6 Scaling effects

As discussed in Section 1.2, to have convergence of the right-hand side of equation (1), the largest absolute eigenvalue of the direct dependency matrix G_{dir} (say β) should be strictly smaller than one (i.e., $\beta < 1$). However, this matrix is unknown in advance and in fact, the whole purpose of having the described model in equation (1) is to compute the direct dependency matrix. In this section, we show that by linearly scaling the un-scaled observed dependency matrix G_{obs}^{us} , it can be guaranteed that the largest absolute eigenvalue of the direct dependency matrix is strictly less than one. In the following, we describe how the linear scaling factor α is chosen where $G_{obs} = \alpha G_{obs}^{us}$.

Suppose $\lambda_{+}^{obs(us)}$ and $\lambda_{-}^{obs(us)}$ are the largest positive and smallest negetive eigenvalues of G_{obs}^{us} . Then, by having

$$\alpha \le \max\Big(\frac{\beta}{(1-\beta)\lambda_{+}^{obs(us)}}, \frac{-\beta}{(1+\beta)\lambda_{-}^{obs(us)}}\Big), \qquad (12)$$

the largest absolute eigenvalue of G_{dir} will be less or equal to $\beta < 1$.

In the following, we show how inequality (12) is derived. Say λ^{dir} and $\lambda^{obs(us)}$ are eigenvalues of direct and unscaled observation matrices. By a similar argument as the one of equation (2), we have

$$\lambda^{dir} = \frac{\lambda^{obs(us)}}{\frac{1}{\alpha} + \lambda^{obs(us)}}.$$
(13)

This function is depicted in Figure S1-b. We consider three cases: $\lambda^{obs(us)} \geq 0$, $-1/\alpha < \lambda^{obs} < 0$ and $\lambda^{obs(us)} < -1/\alpha$ (in the case of $\lambda^{obs(us)} = -1/\alpha$, the function is undefined. Hence, we choose α to avoid this case).

• Case 1: when $\lambda^{obs(us)} \ge 0$: In this case, we have

$$\frac{\lambda^{obs(us)}}{\frac{1}{\alpha} + \lambda^{obs(us)}} \le \beta < 1 \Rightarrow \alpha \le \frac{\beta}{(1-\beta)\lambda^{obs(us)}}$$

• Case 2: when $-1/\alpha < \lambda^{obs} < 0$: In this case, we have

$$\frac{-\lambda^{obs(us)}}{\frac{1}{\alpha} + \lambda^{obs(us)}} \le \beta < 1 \Rightarrow \alpha \le \frac{-\beta}{(1+\beta)\lambda^{obs(us)}}$$

• Case 3: when $\lambda^{obs(us)} < -1/\alpha$:

In this case, we have

$$\frac{\lambda^{obs(us)}}{\frac{1}{\alpha} + \lambda^{obs(us)}} \le \beta < 1 \Rightarrow \frac{1}{\alpha} \le \frac{1 - \beta}{\beta} \lambda^{obs(us)} < 0,$$

which is not possible since $\alpha > 0$. Therefore, α should be chosen so that eigenvalues of the unscaled observation matrix do not fall in this regime. In other words, for negative $\lambda^{obs(us)}$, $\alpha < \frac{-1}{\lambda^{obs(us)}}$. This condition trivially holds as a consequence of Case 2.

Putting all three cases together, inequality (12) is derived which guarantees the largest absolute eigenvalue of the direct dependency matrix is less than or equal to $0 < \beta < 1$.

In equation (1), diffusion effects decay exponentially with respect to the indirect path length. For example, second order indirect effects decay proportional to β^2 where β is the largest absolute eigenvalue of G_{dir} . Therefore, smaller β means faster convergence of the right-hand side of equation (1) and in turn faster decay of diffusion effects. In other words, if β is very small, higher order interactions play insignificant roles in observed dependencies since they decay proportionally to β^k where k is the order of indirect interaction. Figure S4 demonstrates the effect of β in the performance of three different applications: inferring gene regulatory network, inferring protein evolutionary constraints and inferring weak and strong social ties. As it is illustrated in this figure, choosing β close to one (i.e., considering higher order indirect interactions) leads to the best performance in all considered network deconvolution applications. Note that, β has higher effects in the gene regulatory networks and protein evolutionary constraint networks applications compared to the co-authorship network application, which means low-order indirect interactions are sufficient to characterize indirect flows in the coauthorship network compared to two other applications. This is mainly due to the special structure that we observe in the co-authorship network: The co-authorship network can be decomposed locally into densely connected subgraphs and thus most transitive edges happen locally within each module, leading to small indirect effects for longer paths, and thus small β values seem sufficient to characterize their indirect flows.

1.7 Computational complexity analysis of network deconvolution

The proposed network deconvolution algorithm based on eigen decomposition requires computations of eigen decomposition of the observed dependency matrix. This step has order of $\mathcal{O}(n^3)$ computational complexity where n is the number of nodes of the network. Further, the gradient descent based algorithm is iterative and can be very costly to use. Therefore, using baseline network deconvolution algorithms for very large networks can be unattractive due to this required computational power. Here, we propose two extensions to network deconvolution for very large networks. First, many of these large networks are usually low rank and have many zero eigenvalues and therefore the eigen decomposition step can be performed efficiently by exploiting this sparsity. Second, network deconvolution can be parallelized by being applied on local subgraphs of the network. In this section, we explain these two algorithms in more detail.

To show the effectiveness of these extensions, we generate simulated scale-free networks with various levels of eigen sparsity and number of nodes. We add some transitive edges to the network probabilistically by only considering two-hop and three-hop indirect paths. Then, we apply ND to discover transitive edges.

1.7.1 Using sparsity in eigen decomposition

In many applications, networks are low rank (many of network eigenvalues are zero or negligible.). In those cases, eigen decomposition can be performed efficiently by only computing non-zero eigenvalues³. Here, we use MAT-LAB's sparse eigen decomposition algorithm^{15,16}. Figure S5-a shows CPU time usage of performing network deconvolution on networks with 1000 to 50000 nodes and various eigen sparsity levels. As it is illustrated in this figure, network deconvolution can be performed efficiently even on very large networks if the network is low rank. Figure S5-b shows the CPU time usage of ND for various network sizes averaged over networks whose eigen sparsity is less than or equal to 100. This figure demonstrates that unlike the baseline network deconvolution algorithm which has order $\mathcal{O}(n^3)$ complexity, the complexity of performing network deconvolution on low rank networks is approximately linear with respect to the number of nodes in the network.

1.7.2 Parallelization of network deconvolution

Parallelization of algorithms that require large computational power is important in practice. Here, we explain how network deconvolution can be performed in parallel for very large networks:

First, network is divided into subgraphs (perhaps overlapping) with certain sizes. For each subgraph, first and second neighbors of nodes that are not in that subgraph are called marginal nodes of that subgraph. To make sure that edges that are close to subgraph borders do not suffer from network partitioning, we add marginal nodes to each subgraph and then perform network deconvolution. Then, we keep direct dependencies only for subgraph edges (and not for marginal ones). To merge results of different subgraphs, if edges are overlapping, we take the mean of their direct dependencies obtained from different subgraphs. Figure S5-c demonstrates the total CPU time usage and also CPU time usage per core (per subgraph) for different subgraph sizes of a network with 5000 nodes.

2 Inferring Gene Regulatory Networks by Network Deconvolution

In this section, we focus on the application of network deconvolution on inferring gene regulatory networks by using high throughput gene expression data. In gene regulatory network inference, network deconvolution *denoises* the network structure; it eliminates indirect information flows over transitive and noisy edges. Network deconvolution can be used either as a stand-alone individual network inference method or to improve published inference algorithms. In this section, we present the considered framework and demonstrate detailed evaluation results of network deconvolution.

2.1 DREAM framework

we use datasets and the framework from the DREAM5 project⁸ to evaluate the performance of the proposed network deconvolution algorithm as a regulatory inference method. Inferring genome-scale transcriptional regulatory networks from gene expression microarray compendia for a prokaryotic model organism (*E. coli*), a eukaryotic model organism (*S. cerevisiae*) and an *in silico* benchmark is considered in our evaluations. Each of these compendia is represented as an expression matrix of genes and their chip measurements. Further, a list of candidate transcription factors for each compendium is provided. These datasets are fully anonymized (Figure S6).

For each network, a list of directed, unsigned edges are predicted, ordered according to confidence scores. The confidence scores are only to order predicted edges and are not used directly in the scoring. To score each of these predicted networks, organism specific gold standards containing the known transcription factor to target gene interactions (true positives) are used as compiled in the DREAM project⁸. These gold standards are not available in the inference part. For the evaluation, all transcription factor-target gene pairs that are not part of the gold standards are considered as negatives, although, as the gold standards are based on incomplete knowledge, they might contain yet unknown true interactions.

Two considered biological networks (*E. coli* and *S.cerevisiae*) have 4297 and 5667 genes where 296 and 183 of them are transcription factors, respectively. The *in silico* network has 1643 genes where 195 of them are transcription factors. For each of these networks, gene expressions are reported under several conditions such as time courses, perturbations, knock-outs, etc.

It has been shown that, preprocessing the data by assigning different weights to various expression data types can enhance the quality of predicted networks⁸. However, we avoid data preprocessing to demonstrate the gain obtained by network deconvolution algorithm alone.

A key assumption in inferring gene regulatory networks from gene expression data is that, mRNA levels of regulators and their targets show some mutual dependencies. Therefore, the dependency content within interactions of a network can be used as a coarse estimate of inference difficulty. *E. coli* and *in silico* networks demonstrate higher information contents over interacting pairs compared to non-interacting ones⁸. This is true regardless whether Pearson's correlation or mutual information is used to measure dependencies. However, in *S.cerevisiae* network, the dependency distribution of interacting and non-interacting pairs are almost identical. This shows it is more difficult to infer the regulatory network from gene expression data. This fact is reflected in low scores of various inference methods for this network.

2.2 Mutual Information

Mutual information is a quantity that measures the mutual dependence of two random variables. For discrete random variables X and Y, it can be defined as:⁴

$$I(X;Y) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x,y) \log\left(\frac{p(x,y)}{p(x)p(y)}\right), \quad (14)$$

where p(x, y) is the joint probability distribution of X and Y over their supports \mathcal{X} and \mathcal{Y} , and p(x) and p(y) are the marginal probability distribution functions of X and Y, respectively.

In the case of continuous variables, we use a mutual information implementation based on B-spline smoothing of order 3 and discretizing into 10 $bins^5$.

2.3 Evaluation functions

We use the same performance evaluation metrics as in the DREAM5 challenge⁸. Network predictions are evaluated as binary classification tasks where edges are predicted to be present or absent. Then, standard performance metrics from machine learning are used: precision-recall (PR) and receiver operating characteristic (ROC) curves. Then, predictions are evaluated statistically and an overall score summarizing the performance over all three networks is computed. Note that, similar to DREAM5 framework, for each network only top 100,000 edges are considered in the performance evaluation. However, to perform network deconvolution, dense input matrices (~10% density) are used.

For each predicted network, the performance was assessed by the area under the ROC curve (AUROC, true positive rate vs. false positive rate), and the area under the precision vs. recall curve (AUPR). Expressions for true positive rate (TPR), false positive rate (FPR), precision and recall as a function of the cutoff (k) in the edge list are as follows⁸:

$$\begin{aligned} recall(k) &= \frac{TP(k)}{P}, \\ precision(k) &= \frac{TP(k)}{TP(k) + FP(k)} = \frac{TP(k)}{k} \end{aligned}$$

where TP(k) and FP(k) are the number of true positives and false positives in the top k predictions, and P is the number of positives in the gold standard. The *true positive rate*, TPR(k), is similar to *recall*. The *false positive rate* is the fraction of negatives incorrectly predicted in the top k predictions:

$$FPR(k) = \frac{FP(k)}{N},$$
(15)

where N is the number of negatives in the gold standard.

AUROC and AUPR values were separately transformed into *p*-values by simulating a null distribution for a large number (25,000) of random networks⁸. Overall *ROC* and *PR* scores are derived by calculating the geometric mean of the network specific *p*-values:

ROC score =
$$\frac{1}{3} \sum_{i=1}^{3} -\log_{10} p_{ROC_i}$$

PR score = $\frac{1}{3} \sum_{i=1}^{3} -\log_{10} p_{PR_i}$

where p_{ROC_i} and p_{PR_i} represent *p*-values of ROC and PR curves for network *i*. Then, an overall score summarizing the performance over all three networks is calculated as the mean of the ROC and PR scores.

2.4 Network motif analysis

We performed a network motif analysis^{10,8} to compare the ability of inference methods to correctly predict feedforward loops and regulatory cascades before and after ND (**Figure 2b**). We used the refined approach introduced by Petri et al., which measures performance for different motif types based on the area under the ROC curve (AUROC)⁹. Briefly, this involves: (1) identifying the subset of edges pertaining to a given motif type in the true network; (2) evaluating the performance for this subset of edges using the AUROC. Note that edges are counted only once for a given motif type, even if they are part of several overlapping motif instances of the same type.

Using this approach, we computed for each network inference method the motif-specific AUROC of regulatory cascades and feed-forward loops. The prediction bias is given by the difference between the overall AUROC of a given inference method (computed over the complete set of edges as described in Methods) and the motif-specific AUROC. **Figure 2b** shows the average prediction biases over the *in silico* and *E. coli* networks (following the DREAM project⁸, we excluded *S.cerevisiae* because the inference methods did not recover sufficient true edges for a meaningful network motif analysis).

3 Inferring Protein Structural Constraints by Network Deconvolution

Application of network deconvolution to pairwise evolutionary correlation matrix reduces transitive noise effects which in turn leads to higher quality contact map predictions. We evaluate the performance of network deconvolution in inferring contact maps for fifteen tested proteins in different folding classes with sizes ranging from 50 to 260 residues¹¹. The list of these proteins are presented in Table S2. Mutual information (MI) and direct informa $tion^{11}$ (DI), a method recently developed for the proteinstructure prediction problem that seeks to explicitly maximally distribute the interaction probabilities across different residue pairs were used to compute co-variation similarities among different residues. We used direct information and alignment methods provided in http://evfold.org. Then, network deconvolution is applied on MI and DI matrices.

We say two amino acid pairs *contact* if their CA atom distance is less than 7 Angstrom. Trend of our results is robust to the choice of threshold (Figure S10). However, choosing a small threshold leads to sparse contact maps over which significant statistical results cannot be inferred, while considering large thresholds may introduce false contacts. We also report the results for both smaller (5 Angstrom, Figure S10-a) and larger (8 Angstrom, Figure S10-b) thresholds. Distant contact maps are the contacts whose amino acid indices differ at least by 4 with two modifications. We exclude contacts in alpha-helices between positions i and i+8 when there are already contacts between positions i and i+4, and between i+4 and i+8. We also exclude contacts between positions i and i+4 to remove trivial contacts due to alpha-helices. More elaborate contact refinement schemes can also be used as well¹¹.

For evaluation, we consider fraction of discovered true positive contacts among top predictions. We also say a discovered contact is *non-redundant* when it cannot be explained by previously discovered ones (i.e., there is no indirect path connecting them over the network). Non-redundant contacts provide important constraints in predicting 3D structure of proteins.

We compare the performance of network deconvolution (ND) against mutual information (MI) and direct information (DI). Figure S7 shows average total and non-redundant contact discovery rates for 14 tested proteins. We excluded protein 1hzx as all methods had very low performance on it possibly due to its poor alignment.

Recovery rate of contacts after ND is consistently higher than MI for each of 15 tested proteins, leading to an av-

erage increase of 90%, 65% and 50% in the top 50, 100 and 250 predictions. ND leads to a small but consistent improvement across the top predictions of direct information (DI). Remarkably, in the strongest-scoring predictions, ranked between 0 and 50 for each method, ND outperforms DI, both independently and as a post-processing step (Figure S7-a and S8). Focusing on discovery of nonredundant contacts, that cannot be explained by previously discovered ones, ND shows an even greater improvement over the performance of MI, by 116%, 95%, and 60% in the top 50, 100, and 250 predictions, respectively (S7b). Considering non-redundant interactions, ND shows improvement of 9%, 8%, and 6% over DI for the top 50, 100, and 250 predictions, respectively (Figure S7-b and Figure S9). We also developed an alternate definition of redundant contacts, using a window-based approach (i.e., predictions that are within a distance window of previously-discovered contacts are considered redundant). In this metric, ND shows consistent performance increase at all cut-offs over MI and DI (Figure S7-c,d).

4 Inferring Weak and Strong Collaboration Ties by Network Deconvolution

Network deconvolution provides a systematic way of distinguishing weak and strong ties in the network. Application of network deconvolution on an input adjacency matrix ranks edges based on their global importance in information spread through the network. Top ranking edges have greater influence on information spread in the network, while deletion of low ranking edges have a minor effect on network information flows.

In this section, we consider the application of network deconvolution in inferring weak and strong ties in a coauthorship network formed by M. Newman¹³. Here, we introduce this network and a gold standard which assigns tie strengths to edges by using additional information about each publication.

The co-authorship network created by M. Newman¹³ is used to evaluate the performance of network deconvolution in distinguishing weak and strong ties. Nodes of this network represent 1589 scientists that work in the network science field. Two scientists are considered connected if they have coauthored at least one paper.

Additional publication details such as the number of papers each pair of scientists has collaborated and the number of coauthors on each of those papers have been used to calculate a *collaboration tie strength* for each pair of scientists as follows¹⁴:

First, each joint publication is weighted inversely accord-

ing to its number of coauthors. The intuition is that, in general two scientists whose names appear on a paper with many other coauthors know one another less well than two who were the sole authors of a paper. Secondly, the strengths of the ties derived from each of the papers written by a particular pair of scientists are added together. The intuition of this part is that, authors who have written many papers together know one another better on average than those who have written few papers together. Obviously, this is only a coarse approximation of calculating collaboration tie strengths.

As explained above, additional information about each publication can be used to calculate tie strengths¹⁴. However, we show that, network deconvolution can describe many of these collaboration tie strengths solely by use of the topology of the network.

We say two individuals have a strong coauthorship tie if their assigned coauthorship strength in the calculated gold standard is equal or greater than 1/2. Around 36% of edges are labeled as strong ties according to this threshold (991 out of 2742 edges). The performance trend does not change by changing this threshold.

Application of network deconvolution to the un-weighted co-authorship network allows edges to be ranked based on their global significances in the context of the network. The obtained global edge ranking is significantly correlated with the co-authorship tie gold standards calculated by use of publication details (R^2 correlation coefficient is around 0.76). Note that, in **Figure 4b**, tie strengths are mapped to be between 0 and 1 by using their edge ranks. In particular, network deconvolution allows recovery of around 77% of strong co-authorship ties solely by use of the network topology. Also, transitive collaborations are recognized as weak ties in both ND and Newman's tie strength weights.

References

- G. Altay and F. Emmert-Streib. Inferring the conservative causal core of gene regulatory networks. *BMC Systems Biology*, 4(1):132, 2010.
- [2] D. Bertsekas. Nonlinear programming. 1999.
- [3] Y. Saad. Numerical methods for large eigenvalue problems. Manchester University Press. 1992.
- [4] T. Cover, J. Thomas, J. Wiley, et al. *Elements of information theory*, volume 6. Wiley Online Library, 1991.
- [5] J. J. Faith, B. Hayete, J. T. Thaden, I. Mogno, J. Wierzbowski, and *et al.* Large-scale mapping and

validation of escherichia coli transcriptional regulation from a compendium of expression profiles. *PLoS Biology*, 5:54–66, 2007.

- [6] T. Hopf, L. Colwell, R. Sheridan, B. Rost, C. Sander, and D. Marks. Three-dimensional structures of membrane proteins from genomic sequencing. *Cell*, 2012.
- [7] R. Horn and C. Johnson. *Matrix analysis*. Cambridge Univ Pr, 1990.
- [8] D. Marbach, et al. Wisdom of crowds for robust gene network inference. Nature Methods, 2012.
- [9] T. Petri, S. Altmann, L. Geistlinger, R. Zimmer, and K. Robert. Inference of eukaryote transcription regulatory networks. *in preparation*, 2012.
- [10] D. Marbach, et al. Revealing strengths and weaknesses of methods for gene network inference. Proceedings of the National Academy of Sciences of the United States of America, 107, 2010.
- [11] D. Marks, L. Colwell, R. Sheridan, T. Hopf, A. Pagnani, R. Zecchina, and C. Sander. Protein 3d structure computed from evolutionary sequence variation. *PloS one*, 6(12):e28766, 2011.
- [12] P. Meyer, D. Marbach, S. Roy, and M. Kellis. Information-theoretic inference of gene networks using backward elimination, 2010.
- [13] M. Newman. Finding community structure in networks using the eigenvectors of matrices. *Physical review E*, 74(3):036104, 2006.
- M. Newman et al. Scientific collaboration networks.
 ii. shortest paths, weighted networks, and centrality. *Physical review-series E-*, 64(1; part 2):16132–16132, 2001.
- [15] R.B. Lehoucq and D.C. Sorensen Deflation Techniques for an Implicitly Re-Started Arnoldi Iteration SIAM J. Matrix Analysis and Applications, Vol. 17, 1996, pp. 789Ũ821.
- [16] D.C. Sorensen Implicit Application of Polynomial Filters in a k-Step Arnoldi Method SIAM J. Matrix Analysis and Applications, Vol. 13, 1992, pp.357Ũ385.



Figure S1: Effects of applying network deconvolution on network eigenvalues

Panel (a) shows an example of how eigenvalues of the observed matrix change by application of network deconvolution. In general, ND decreases magnitudes of high positive eigenvalues that were inflated due to transitivity effects over the network. Panel (b) shows the network deconvolution nonlinear mapping function that is applied to eigenvalues of the unscaled observed dependency matrix to compute the eigenvalues of the direct dependency matrix.



Figure S2: Convergence of network deconvolution for non-decomposable matrices This figure illustrates the convergence of Algorithm 3 over a random asymmetric network with 50 nodes and 490 edges. As the number of iterations increases, the energy of the error term goes to zero.



Figure S3: Robustness of network deconvolution against noise

This figures illustrates the performance of network deconvolution as a gene regulatory network inference method over noisy data. Y-axis represents the inference score (combination of AUPR and AUROC scores) as explained in Section 2.3. However, here, to observe noise effects completely, unlike the DREAM framework scoring (see Section 2.3), we do not threshold inferred networks. *E. coli* (E) and *in silico* (I) networks are considered in this setup where an artificial Gaussian noise with a variance of a fraction of the experiment variance is added to each experiment in gene expression datasets. As it can be seen in this figure, by increasing the noise power, the overall score decreases smoothly which indicates the robustness of network deconvolution against noise. The error bars represent standard deviations of evaluation points.



Figure S4: Linear scaling effects on network deconvolution performance

In the first step of network deconvolution framework, a linear scaling function is applied on the unscaled observed dependency matrix to have the largest absolute eigenvalue of the direct dependency matrix less than or equal to β . Smaller β means faster convergence of the right-hand side of equation (1) and in turn faster decay of diffusion effects. Therefore, in the small β case, higher order indirect interactions play smaller roles in observed dependencies. This figure illustrates the effect of β in the performance of network deconvolution in various considered applications: Panel (a) demonstrates the overall score of network deconvolution as a gene regulatory network inference method, over three considered networks E. coli, in silico and S. cerevisiae. This overall score is computed by combining AUPR and AUROC scores as explained in Section 2.3. For more details, see Section 2. Panel (b) demonstrates the effect of β in the performance of network deconvolution in inferring protein evolutionary constraints. Y-axis shows the overall contact map prediction rate of top 100 predictions over all 15 tested proteins as explained in Section 3. For more details, see Section 3. Panel (c) illustrates the effect of β in the performance of network deconvolution in inferring weak and string social ties. Y-axis shows the strong tie prediction ratio over the considered co-authorship network as explained in Section 4. As it is illustrated in this figure, choosing β close to one (i.e., considering higher order indirect interactions) leads to best performance in all considered network deconvolution applications. Further, note that despite the social network application, considering higher order indirect interactions is important in gene regulatory inference and protein structural constraint inference applications.





Panel (a) shows CPU time usage of performing network deconvolution on networks with 1000 to 50000 nodes and various eigen sparsity levels. As it is illustrated in this figure, network deconvolution can be performed efficiently even on very large networks if the network is low rank. Panel (b) shows the CPU time usage of ND for various network sizes averaged over networks whose eigen sparsity is less than or equal to 100. This figure demonstrates that unlike the baseline network deconvolution algorithm which has a complexity of order $\mathcal{O}(n^3)$, complexity of performing network deconvolution on low rank networks is approximately linear with respect to the number of nodes in the network. (c) Parallelization of network deconvolution is important to apply it effectively on very large networks. This panel demonstrates the total CPU time usage and also CPU time usage per core (per subgraph) for different subgraph sizes of a network with 5000 nodes.





DREAM5 participants were solicited to infer gene regulatory interactions from three different expression compendia. For each network, a list of directed, unsigned edges are predicted, ordered according to confidence scores. To score each of these predicted networks, organism specific gold standards containing the known transcription factor to target gene interactions (true positives) are used as compiled in the DREAM project⁸. These gold standards are not available in the inference part. For the evaluation, all transcription factor-target gene pairs that are not part of the gold standards are considered as negatives.





(a) Recovery of truly interacting amino acid residues with increasing numbers of top-scoring pairs for mutual information (MI), direct information (DI), and network deconvolution (ND). Recovery after ND is consistently higher than MI for each of tested proteins leading to an average increase of 90%, 65% and 50% in the top 50, 100 and 250 predictions. ND leads to a small but consistent improvement across the top predictions of direct information (DI). In the strongest-scoring predictions, ranked between 0 and 50 for each method, ND outperforms DI, both independently and as a post-processing step. (b) Discovery of non-redundant contacts, that cannot be explained by previously discovered ones. ND shows improvement over both MI and DI specially in top predictions. (c),(d) We show here an alternate definition of redundant contacts, using a window-based approach (predictions that are within a distance window of previously-discovered contacts are considered redundant). In this metric, ND shows consistent performance increase at all cut-offs.





Contact recovery rate after ND is consistently higher than MI for each of 15 tested proteins, leading to an average increase of 90%, 65% and 50% in the top 50, 100 and 250 predictions. ND leads to a small but consistent improvement across the top predictions of direct information (DI). Remarkably, in the strongest-scoring predictions, ranked between 0 and 50 for each method, ND outperforms DI, both independently and as a post-processing step. We excluded protein 1hzx (highlighted in the figure) from overall performance comparisons as all methods had very low performance on it possibly due to its poor alignment.



top predictions



Focusing on non-redundant contacts, that cannot be explained by previously discovered ones, ND shows a significant improvement over the performance of MI by 116%, 95%, and 60% in the top 50, 100, and 250 predictions, respectively. The improvement of ND over DI is even more pronounced when considering non-redundant interactions, with ND showing improvement of 9%, 8%, and 6% for the top 50, 100, and 250 predictions, respectively.



Figure S10: Total and non-redundant discovery rate for different thresholds of contact proximity Performance trend of ND is robust to the choice of proximity threshold to define contacts. However, choosing a small threshold leads to sparse contact maps over which significant statistical results cannot be inferred, while considering large thresholds may introduce false contacts. We choose threshold 7 Angstrom and also report the results for both smaller (5 Angstrom, part a) and larger (8 Angstrom, part b) thresholds.

Network name	# of nodes	Network density	Ave. # of neighbors	Clustering coeff.	Network diameter	Network centrality
in silico	1643	0.003	5.107	0.185	10	0.021
E. coli	4297	0.004	3.802	0.19	13	0.252
S. cerevisiae	5667	0.002	3.947	0.071	9	0.108
Co-authorship	1589	0.003	3.754	0.694	17	0.021

 Table S1: General properties of considered networks.

This table presents general properties of in silico, E. coli, S. cerevisiae and coauthorship networks.

Protein PDB ID	Protein name	Protein PDB ID	Protein name	Protein PDB ID	Protein name
1hzx	OPSD BOVIN	1bkr	SPTB2 HUMAN	1odd	OMPR ECOLI
3tgi	TRY2 RAT	2it6	A8MVQ9 HUMAN	1g2e	ELAV4 HUMAN
5p21	RASH HUMAN	1rqm	THIO ALIAC	1wvn	PCBP1 HUMAN
1f21	RNH ECOLI	2072	CADH1 HUMAN	5pti	BPT1 BOVIN
1e6k	CHEY ECOLI	1r9h	O45418 CAEEL	2hda	YES HUMAN

Table S2: Tested protein names and PDB IDsThis table presents tested protein names and their PDB IDs.