

Submitted to the *Annals of Applied Statistics*  
arXiv: [stat.AP/0907.0000](https://arxiv.org/abs/1409.0700)

## REFINING CELLULAR PATHWAY MODELS USING AN ENSEMBLE OF HETEROGENEOUS DATA SOURCES\*

BY ALEXANDER M. FRANKS<sup>§</sup>, FLORIAN MARKOWETZ\*<sup>¶</sup>  
AND EDOARDO M. AIROLDI\*<sup>§</sup>

<sup>§</sup>*Harvard University and* <sup>¶</sup>*University of Cambridge*

Improving current models and hypotheses of cellular pathways is one of the major challenges of systems biology and functional genomics. There is a need for methods to build on established expert knowledge and reconcile it with results of high-throughput studies. Moreover, the available data sources are heterogeneous and need to be combined in a way specific for the part of the pathway in which they are most informative. Here, we present a compartment specific strategy to integrate edge, node and path data for the refinement of a network hypothesis. Specifically, we use a local-move Gibbs sampler for refining pathway hypotheses from a compendium of heterogeneous data sources, including novel methodology for integrating protein attributes. We demonstrate the utility of this approach in a case study of the pheromone response MAPK pathway in the yeast *S. cerevisiae*.

**1. Introduction.** Cellular mechanisms are driven by interactions between DNA, RNA, and proteins working together in cellular pathways. However, the current knowledge of information flow in the cell is still very incomplete (Kirouac et al., 2012). Even in well established signaling pathways studied for decades in model organisms, newer approaches can discover novel components (Müller et al., 2005) or cross-talk with other pathways (McClean et al., 2007; Vaga et al., 2014). In cancer, finding pathways underlying disease development can lead to new drug targets (Balbin et al., 2013). This makes the dissection of cellular pathways one of the major challenges of systems biology and functional genomics. We can represent cellular pathways using a network, in which edges represent conditional dependence between molecule abundances and non-edges represent conditional independence. Thus, the challenge is to infer or refine a hypothesis about the edges in pathway model.

---

\*This work was supported, in part, by NIH grant R01 GM-096193, NSF CAREER grant IIS-1149662, and by MURI award W911NF-11-1-0036 to Harvard University. EMA is an Alfred P. Sloan Research Fellow and a Shutzer Fellow at the Radcliffe Institute for Advanced Studies. FM acknowledges support from the University of Cambridge, Cancer Research UK (C14303/A17197), and Hutchison Whampoa Limited. FM and EMA contributed equally to this work.

*Signalling pathways.* In this paper, we focus on signalling pathways, which are of particular importance because they encode how cell to reacts to external stimuli (Alberts et al., 2002). Starting from receptor proteins in the *cell membrane* these pathways traverse the *cytoplasm* by relaying the signal from one protein to the next, often by phosphorylation in so called MAPK cascades. At the end of these cascades lie *transcription factors*, which are specialized proteins that move from the cytoplasm into the *cell nucleus*, bind there to DNA and regulate gene expression as a response to the external stimulus. Signaling pathways thus traverse and connect the major compartments of the cell: the membrane, cytoplasm and the nucleus.

*Inferring signalling pathways from data.* One of the main obstacles to utilize high-throughput data in refining known pathway models is the gap between the relatively unbiased and hypothesis-free nature of generating genome-scale datasets and the need for very focused, hypothesis-driven research to test biological models in small or medium scale experiments (Hibbs et al., 2008). While researchers in computational biology usually start with a collection of data and reconstruct pathways from it, experimental biologists often start with a specific network hypothesis in mind and try to reconcile it with the evidence from high-throughput screens.

*Our approach.* Here, we contribute to bridging this gap by introducing a comprehensive data integration strategy to refine a given network hypothesis. Our approach is characterized by three key features, which set it apart from previous approaches: First, we start with a *specific pathway model* (represented by a network) and assess how well it is supported in a collection of complementary data sets. These data sets are heterogeneous and informative for distinct cellular locations. Second, we exploit this fact by introducing a *compartment-specific* probabilistic model, which distinguishes different cellular locations (the membrane, cytoplasm and nucleus) and where data types are only used for reconstructing the parts of the network they are informative about. Third, we explicitly include *node properties* in our model. This allows us to use data on the properties of the molecules like protein phosphorylation states or protein domains, which have so far been under-utilized for pathway structure learning (Ryan et al., 2013).

In this paper we show that our modeling approach can assist experimentalists in planning future studies by assessing which parts of a biological model are not well supported by data, and by proposing testable extensions and refinements of a given pathway hypotheses. We demonstrate the power of our approach in a case study in the yeast *S. cerevisiae*.

*Related work.* Pathway reconstruction is a well established field in compu-

tational biology and statistics (Hyduke and Palsson, 2010; Markowitz and Spang, 2007). Several features distinguish our pathway refinement methodology from existing network reconstruction methods.

Comprehensive data integration strategies on large data collections were shown to be very successful in predicting protein function and interactions (Guan et al., 2012; Llewellyn and Eisenberg, 2008; Guan et al., 2008; Myers et al., 2005). These methods are very helpful for describing the global landscape of protein function, but offer less insight into individual molecular mechanisms and pathways. Our approach differs from methods to refine pathway hypotheses from expression profiles of down-stream regulated genes (Gat-Viks and Shamir, 2007), because we integrate heterogeneous data sources in a compartment-specific way.

We also differ from previous research on de-novo pathway reconstruction. These methods can be classified by how they use information about edges, paths and nodes in the pathway diagram for structure learning.

- **Edges:** Most approaches incorporate evidence for individual edges in the network using correlation measures (Mulder et al., 2012; Wang et al., 2012; Li et al., 2013) or higher-order graphical models (Schäfer and Strimmer, 2005a; Friedman, 2004; Segal et al., 2003), sometimes integrating additional data sources into the model (Bernard and Hartemink, 2005; Werhli and Husmeier, 2007; Balbin et al., 2013; Gitter et al., 2013; Nariai et al., 2004; Segal et al., 2003).
- **Paths:** Cause-effect relationships indicating paths from perturbed genes to observed effects are exploited in methods like SPINE (Ourfali et al., 2007), physical network models (Yeang et al., 2005), nested effects models (Wang et al., 2014; Markowitz et al., 2007; Tresch and Markowitz, 2008; Fröhlich et al., 2007, 2008) and others (Lo et al., 2012; Yip et al., 2010), with applications including DNA damage repair (Workman et al., 2006) and cancer signalling (Knapp and Kaderali, 2013; Stelniec-Klotz et al., 2012).
- **Nodes:** Features of individual proteins or genes provide data for nodes and have been found useful for predicting that a protein contributes to a pathway (Hahne et al., 2008; Fröhlich et al., 2008) but have so far been under-utilized in reconstructing pathway structure (Ryan et al., 2013).

Our method differs from existing methods in several important aspects: First of all, we are the first to integrate data about *edges* and *paths* as well as *nodes* in the pathway diagram. Additionally, in contrast to de-novo network reconstruction we start with a hypothesis network and identify which hypothesized edges are supported by the data. We also differ from other

methods which evaluate formal one and two sample network hypothesis tests (Yates and Mukhopadhyay, 2013). Our goal is not to explicitly determine whether our initial hypothesis is “correct”—on the contrary we assume a priori that any initial hypothesis can be further refined and improved upon. We provide a list of edge probabilities that can assist experimentalists in their future studies. We assess which parts of an existing biological model are not well supported by a data as well as suggesting new edges which are supported by the data but which are not part of the original hypothesis.

*Overview.* We describe a compartment-specific probabilistic graphical model for posterior inference on cellular pathways in *section 2*, which can be used to extend and refine a given biological model and predict novel parts of the pathway graph. Our model comprehensively integrates the three general types of data on edges, paths, and nodes. We demonstrate the utility of our methods in a case study in *S. Cerevisiae* (*section 3*) by first exploring how informative different data sources are individually (*section 3.1*) and then evaluating results of posterior draws using both full data and leave-one-out data (*section 2*).

**2. An integrative model of a cellular pathway.** Given a set of a gene products, i.e., putative pathway members, we infer an undirected network model using a local-move Gibbs sampler. The network model, is defined in terms of  $N$  nodes and the edges between these pairs of nodes,  $(n, m)$ . The edges are encoded by a binary random variable,  $X_{nm}$ . The collection of edge-specific random variables defines the adjacency matrix,  $\mathbf{X}$ , of the pathway model.

*Parameter estimation and posterior inference.* The adjacency matrix  $\mathbf{X}$  corresponding to the pathway model is latent since we cannot directly observe the edges, though we have strong prior belief about many edges. Thus, the primary goal of our analysis is to do posterior inference on the adjacency matrix,  $\mathbf{X}$ , from a collection of  $\mathcal{K}$  data sets,  $Y_{1:\mathcal{K}}$  and an initial pathway hypothesis. Although we treat  $\mathbf{X}$  as latent, we differ from de-novo pathway reconstruction by incorporating this informative hypothesis pathway which we use to train the models for data sets  $Y_{1:\mathcal{K}}$  (see Section 3).

By Bayes rule, the posterior distribution on a pathway model,

$$(2.1) \quad P(X | Y_{1:\mathcal{K}}, \Theta) \propto P(X | \Theta) \cdot P(Y_{1:\mathcal{K}} | \mathbf{X}, \Theta),$$

is proportional to the prior distribution on the pathway with the likelihood of the data. Here,  $\Theta$  is a collection of parameters for the data models introduced below.

We use a local Gibbs sampling strategy to sample pathway models from posterior distribution in Equation 2.1. The sampler explores the space of pathway models by adding or removing edges in turn, one at a time. Specifically, the edge  $X_{nm}$  between gene products  $(n, m)$  is sampled according to a Bernoulli distribution, with probability of success

$$(2.2) \quad P(X_{nm} \mid X_{(-nm)}, Y_{1:\mathcal{K}}, \Theta),$$

where  $X_{(-nm)}$  represents the set of edges without  $X_{nm}$ .

2.1. *Context-specific data contributions through a compartment map.* We use five complementary data types: physical binding of protein pairs (including yeast-two hybrid, mass spectrometry, and literature-curated data), transcription factor-DNA binding assays, gene knockout data, gene co-expression data, and node information (including protein domains and differential phosphorylation arrays).

Importantly, different data sets can be very informative in specific cellular locations while completely uninformative in others. Thus, before we define the data likelihoods in section 2.2, it is essential to exploit this fact in our model.

To instantiate the notion that different data are informative in different cellular locations, we introduce an additional modeling element: the compartment map, which contains three conceptual pathway compartments directly based on the organisation of the cell (Alberts et al., 2002): First, the *cell membrane*, where receptor proteins sense signals from outside the cell; second, the *cytoplasm*, where protein cascades relay these signals to transcription factor proteins that enter the third compartment, the *nucleus*, to regulate the activity of target genes. The compartment map,  $\mathcal{C}$ , is a  $5 \times 3$  binary matrix that associates the three pathway compartments with the five data types to indicate which data type is informative about molecular interactions in which compartments (see Table 1).

In particular, each data set is described by a pair  $D_k = (Y_k, T_k)$ , where  $Y_k$  denotes the collection of measurements, and  $T_k$  is five-level factor that denotes the data type (and indexes the relevant row of  $\mathcal{C}$ ). We can now revise the form of the conditional distributions in Equation 2.2,

$$(2.3) \quad P(X_{nm} \mid X_{(-nm)}, D_{1:\mathcal{K}}, \mathcal{C}, \Theta) = \frac{\mathcal{L}(X_{nm} = 1, X_{(-nm)} \mid D_{1:\mathcal{K}}, \mathcal{C}, \Theta)}{\mathcal{L}(X_{nm} = 1, X_{(-nm)} \mid D_{1:\mathcal{K}}, \mathcal{C}, \Theta) + \mathcal{L}(X_{nm} = 0, X_{(-nm)} \mid D_{1:\mathcal{K}}, \mathcal{C}, \Theta)}$$

Overloading notation, we let  $\mathcal{C}_t(n, m)$  be an indicator reflecting whether data type  $t$  is informative for the protein pair  $(n, m)$ , based on the compartment map and the localizations of proteins  $n$  and  $m$ . This leads to the following likelihood specification:

$$\begin{aligned}
 (2.4) \quad \mathcal{L}(X_{nm}, X_{(-nm)} \mid D_{1:\mathcal{K}}, \mathcal{C}, \Theta) &\propto \\
 &= \prod_k^{\mathcal{K}} \left[ P(Y_k \mid X_{nm}, X_{(-nm)}, T_k = t, \Theta)^{\mathcal{C}_t(n, m)} \right. \\
 (2.5) \quad &\left. \times P(Y_k \mid X_{(-nm)}, T_k = t, \Theta)^{1 - \mathcal{C}_t(n, m)} \right]
 \end{aligned}$$

where the role of the indicator is to discard data collections from data types that are expected to carry little information about the protein pair of interest, according to information in  $\mathcal{C}$ . That is, for any pair  $(n, m)$ ,  $\mathcal{C}_t(n, m) = 0$  implies data set  $Y_k$  is conditionally independent of edge  $(n, m)$  given the rest of the pathway. In this case, the data in  $Y_k$  have no effect on the conditional posterior probability of  $X_{nm}$ .

In Algorithm 1 we outline the steps of the local-move Gibbs sampler. First, we use the initial pathway hypothesis to learn model parameters for the likelihoods described in Section 2.2. These parameters are learned from the hypothesis pathway or another held-out training pathway. For instance, for pheromone pathway inference (Section 3) we can infer these parameters using the hypothesis pheromone pathway or the osmolarity, hypotonic or starvation sub-pathways (Section 3.3). A summary of all data model parameters can be found in Table 1.

After inferring these data parameters, we proceed with the main pathway refinement algorithm. For each pair of vertices in the network (in a randomly chosen order), we sample the presence or absence of an edge from the conditional distribution, given all other edges. As described above, the conditional distribution is based only on the informative data types for the proposed vertices which are determined by the compartment map and cellular locations of the relevant genes.

*2.2. Likelihoods for high-throughput data on edges, paths and nodes.* Data of different types need to be modeled differently. We focus on modeling five main data types: protein interaction data, protein-DNA binding data, gene co-expression data, gene perturbation data, and node attribute data (differential phosphorylation and protein domains). Below, we describe the likelihood functions corresponding to the main data types of interest and methods of inference. For all data models, we use the hypothesis pathway to learn the relevant parameters.

Pathway inference via Gibbs sampling

- 1. Infer model parameters using initial pathway hypothesis**
- 2. Initialize  $X$  to the pathway hypothesis.**

for desired number of samples do

- for  $n, m$  in  $1:N$  do
  - 3. Identify informative data types,  $t \in T$ , using  $C_t(n, m)$**
  - 4. Compute  $\mathcal{L}(X_{nm} = 1, X_{(-nm)}|\cdot)$  and  $\mathcal{L}(X_{nm} = 0, X_{(-nm)}|\cdot)$  using Equation 2.4.**
  - 5. Accept the pathway with  $X_{nm} = 1$  according to Equation 2.3**

**Algorithm 1:** Local-move Gibbs sampler

*Likelihood for protein interaction data.* Here, we consider a single data set  $Y_{N \times N}$  aimed at measuring physical protein binding events (PPI). We reduce the likelihood of the data,  $Y$ , to a function the false positive and false negative rates,  $\alpha$  and  $\beta$ . Given the pathway,  $X$ , we evaluate

$$(2.6) \quad \mathcal{L}_{ppi}(Y | X, \alpha, \beta) = \alpha^{S_{10}}(1 - \alpha)^{S_{11}}\beta^{S_{01}}(1 - \beta)^{S_{00}},$$

where  $S_{xy}$  counts the number of edges for which  $X_{nm} = x$  and  $Y_{nm} = y$ . For instance,  $S_{10}$  is the number of false positives. We estimate  $\alpha$  and  $\beta$  as the maximum likelihood estimates of the appropriate binomial likelihood, e.g.  $\hat{\alpha} = \frac{S_{10}}{S_{10} + S_{11}}$  and  $\hat{\beta} = \frac{S_{01}}{S_{01} + S_{00}}$  where  $S$  can come from the target hypothesis pathway or a different training network.

*Likelihood for protein-DNA binding data.* Here, we consider a single data set  $Y_{N \times M}$  aimed at measuring transcription factor-DNA binding events (TF) of  $N$  genes on  $M < N$  transcription factors. Rather than hybridization levels (for ChIP-chip) or peaks (for ChIP-seq), we model the  $p$ -values corresponding to binding events, which makes our model independent of the technology used to detect the binding event. We develop a mixture model for the  $p$ -values, directly. Given the pathway,  $X$ , we expect to see a small  $p$ -value for protein  $n$  binding nucleotide sequence  $m$  whenever the edge  $X_{nm}$  is present. On the contrary, the  $p$ -values are uniformly distributed under the null hypothesis of no binding events,  $X_{nm} = 0$ . We evaluate

$$(2.7) \quad \mathcal{L}_{tf}(Y | X, \gamma) = \prod_{n,m} [ \text{Uniform}(Y_{nm}) \cdot \mathbb{1}(X_{nm} = 0) + \text{Beta}(Y_{nm} | \gamma, 1) \cdot \mathbb{1}(X_{nm} = 1) ],$$

where  $0 < Y_{nm} < 1$  ( $p$ -value), and  $0 < \gamma < 1$ . See a related beta-uniform mixture model introduced by [Pounds and Morris \(2003\)](#) in the context of

multiple testing for differential expression. For pathway refinement, we take  $\gamma$  to be the maximum likelihood estimate derived from the set of p-values corresponding to edges in the training pathway.

*Likelihood for knock-out data.* Here we consider a data set  $Y_{N \times M}$  with  $M < N$  knockouts, where  $Y_{mn}$  is the log-two-fold change in expression of gene  $n$ , when gene  $s$  is knocked out. Let  $Z_{mn}$  be a binary variable representing the existence of a directed path from gene  $n$  to gene  $m$ , *through a transcription factor*. While we consider the set of undirected pathway models, we temporarily impute directionality using the fact that the cellular signal should flow from the cytoplasm to the nucleus. We model the knockout data as a mixture of normals:

$$(2.8) \quad \mathcal{L}_{ko}(Y | X, \sigma_0, \sigma_1) = \\ = \prod_{n,m} \text{Normal}(Y|0, \sigma_1) \mathbb{1}[Z_{mn}] + \text{Normal}(Y|0, \sigma_0) \mathbb{1}(1 - Z_{mn})$$

The standard deviations for change in expression are represented by  $\sigma_0$  (when there is no path between the knockout and a target) and  $\sigma_1$  (there is a path). Empirically  $\sigma_1 > \sigma_0$  since there is generally a larger change in expression of a node,  $n$ , for knockout  $m$  when  $n$  and  $m$  are connected in the pathway. We take  $\sigma_1$  to be the maximum likelihood estimate based on the set of log-two-fold changes for which there is a direct pathway between the knockout and target in the hypothesis / training pathway. Similarly, we take  $\sigma_2$  to be the maximum likelihood estimate based on the set data for which there is no path between knockout and target.

*Likelihood for gene co-expression data.* Here, we consider a single data set  $Y_{N \times N}$  aimed at measuring gene expression. Rather than hybridization levels (for microarrays) or the number of reads (for mRNA sequencing), we model correlations among the profiles of pairs of genes, which again makes our model independent of the details of the measurement technology. We develop a mixture model for the correlations, directly. Given the pathway,  $X$ , we expect to see correlation between the expression profiles of two genes whenever they are co-regulated. Similarly to [Schäfer and Strimmer \(2005b\)](#), we use a mixture model for the distribution of the sample correlation coefficient  $\hat{\rho} = y$  of the form

$$(2.9) \quad \mathcal{L}_{expr}(Y | X, \delta, \kappa) = \prod_{n < m} [ P_0(Y_{nm} | \kappa) \cdot \mathbb{1}(X_{nm} = 0) + \\ P_1(Y_{nm} | \delta, 1) \cdot \mathbb{1}(X_{nm} = 1) ]$$



When  $X_{nm} = 0$ , we expect the two gene profiles to be uncorrelated. Differently from Schäfer and Strimmer (2005b), however, we chose a distribution that puts more emphasis on higher correlation if we see an edge in the model,  $X_{nm} = 1$ , using a one-parameter beta distribution,

$$(2.10) \quad P_1(y|\delta) = \text{Beta}(y | \delta, 1).$$

As in the model for protein-DNA binding data, we estimate  $\delta$  using maximum likelihood on the set of gene pairs which share a transcription factor in the hypothesis / training pathway.

*Likelihood for node attributes data.* Here, we consider a single data set  $Y_N$  that lists node-specific attributes such as protein domains from PFAM (Punta et al., 2012) and SMART (Schultz et al., 1998; Letunic et al., 2012) databases, and differential phosphorylation data (Gruhler et al., 2005). We develop novel techniques to model protein attributes. Specifically, we model the likelihood of an attribute conditional on the given pathway  $\mathbf{X}$ . We term our models for node attributes “relation regression.” For differential phosphorylation data,  $Y_{N \times 1}$ ,

$$(2.11) \quad \mathcal{L}_{node}(Y | X, \lambda, \sigma) = \prod_n \text{Normal} \left( Y_n \mid \lambda_0 + \lambda_1 \frac{\sum_{m \neq n} Y_m \mathbb{1}(X_{nm} = 1)}{\sum_{m \neq n} \mathbb{1}(X_{nm} = 1)}, \sigma_N^2 \right)$$

In other words, the differential phosphorylation,  $Y_n$ , is assumed to be linearly related to the mean differential phosphorylation of the neighbors of node  $n$ . Similarly, for the protein domain data,  $D_{N \times S}$ , we use an auto-logistic regression to model the data. Specifically, for  $D_{ns}$ , a binary variable indicating the presence of domain  $s$  in protein  $n$ ,

$$(2.12) \quad \mathcal{L}_{node}(D | X, \lambda) = \prod_{ns} P_{ns}^{D_{ns}} (1 - P_{ns})^{(1 - D_{ns})}$$

where

$$P_{ns} = \text{logit}^{-1} \left( \lambda_0 + \sum_j \lambda_j \mathbb{1} \left[ \sum_{m \neq n} D_{mj} \mathbb{1}(X_{nm} = 1) > 0 \right] \right)$$

Here,  $\text{logit}(P_{ns})$  is linearly related to the presence of domains in neighboring genes. In both the normal and logistic regressions, we again fit the parameters  $\vec{\lambda}$ , using training / initial hypothesis pathway. In the logistic model, we

Data type	Parameters
protein interaction	$\alpha, \beta$
protein-DNA binding	$\gamma$
gene knock-out	$\sigma_0, \sigma_1$
gene co-expression	$\kappa, \delta$
node attributes	$\lambda, \sigma_N^2$

TABLE 1

*List of learned parameters for high-throughput data. Prior to pathway refinement, we first infer all parameters using the hypothesis pathway or a distinct “training pathway”.*

use a weakly-informative Cauchy prior for the coefficients (Gelman, 2008). This controls for any overfitting and separation problems.

*Prior distribution on the space of pathway models.* In this study our focus lies on assessing the extent to which the data support a pathway model  $X$ . We choose a block model prior  $P(X)$  over binary matrices of size  $N \times N$  with edge density fixed by compartment. In general, any informative prior distribution on graphs could be used here to encode biological knowledge (Isci et al., 2013; Mukherjee and Speed, 2008).

### 3. Case study: Pheromone Response Pathway in *S. cerevisiae*.

To demonstrate the efficacy of our approach, we examine the pheromone response MAPK pathway in the yeast *S. cerevisiae*. It offers the opportunity to combine a large collection of datasets with a solid understanding of the pathway structure. The pheromone pathway is the subject of intense research efforts in computational biology as well as experimental biology (Hara et al., 2012; Scott et al., 2006; Kofahl and Klipp, 2004) and shows cross-talk to other MAPK pathways (Nagiec and Dohlman, 2012; McClean et al., 2007; Gat-Viks and Shamir, 2007).

*Initial pathway construction.* To start our analysis in a way relevant to refining and extending existing knowledge of signaling pathways, we extracted a model of the pheromone response pathway from the summary of MAPK pathways (sce04010) in the database KEGG (Kanehisa and Goto, 2000) and combined it with known transcription factor (TF) targets from two independent studies (Simon et al., 2001; Ren et al., 2000).

We split the pathway into three parts: the *membrane* compartment containing the receptor proteins, the *cytoplasm* compartment containing the MAPK cascade to activate the transcription factors (TF), and the *nuclear* compartment containing the TFs and their targets. Figure 1-A depicts the pathway hypothesis. Proteins mediating between two compartments (like

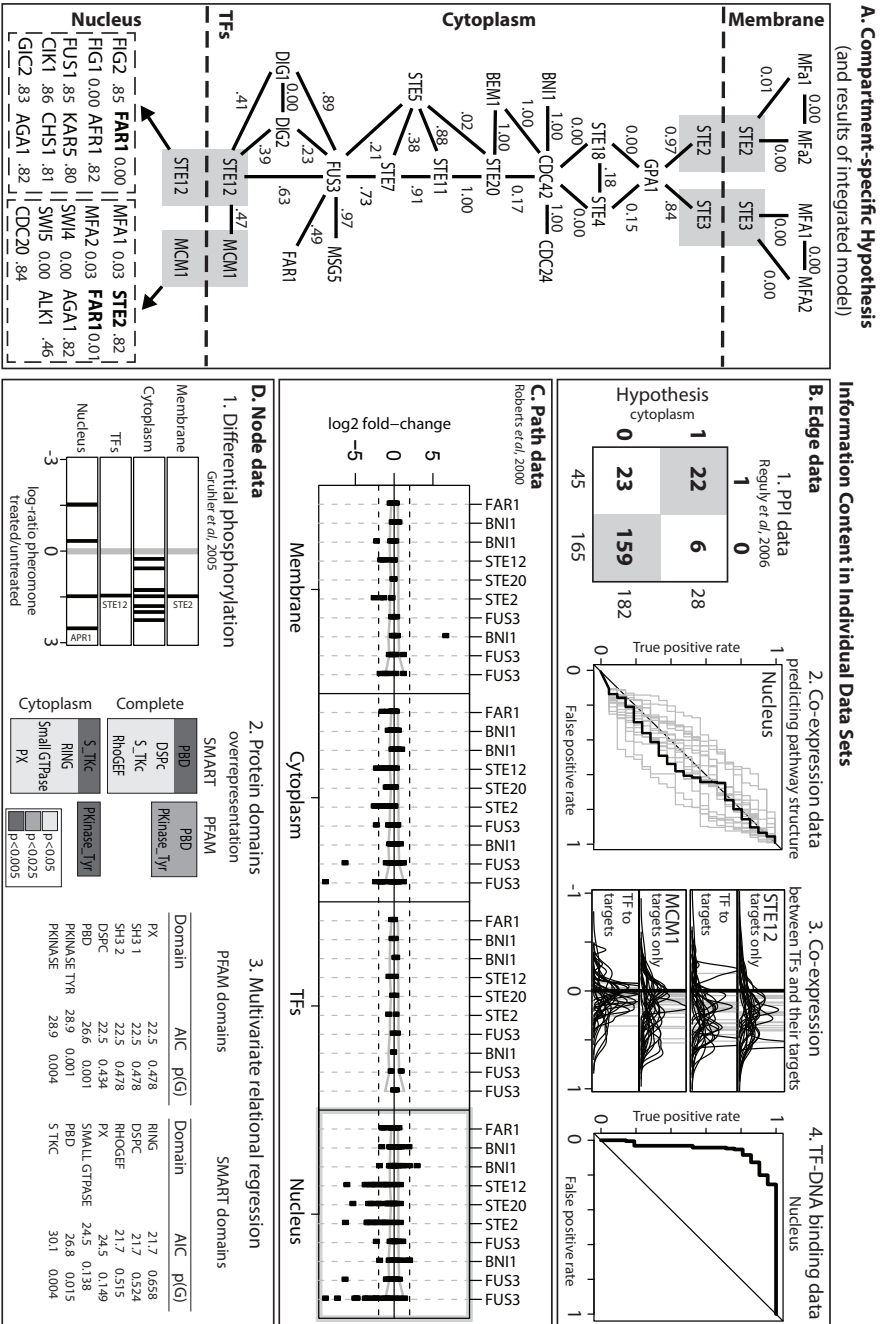


Fig 1: Compartment-specific pathway hypothesis, posterior probabilities, and evaluation of support in the data. **A.** Pathway hypothesis and posterior edge probabilities for the Yeast pheromone response pathway. Dashed lines delineate compartments and grey boxes show pathway members active in more than one compartment. Panels B-D summarize to which extent the hypothesis is reflected in individual data types. **B.** Edge data: (1) protein-protein interactions in the cytoplasm, (2) ROC curve using gene co-expression in the nucleus, (3) co-expression of TFs with their targets in the cytoplasm, and (4) ROC curve using TF binding data in the nucleus. **C.** Cause-effect data shows different transcriptional effects in the four compartments, with nuclear effects being most prominent. **D.** Node data: (1) Differential phosphorylation, (2) Overrepresentation of protein-domains in different compartments, (3) goodness-of-fit of auto-logistic models on protein domains from PFAM and SMART.

TFs) are contained in two sub-graphs and marked by grey boxes. TF targets that are also members of other compartments are indicated in bold.

3.1. *Exploratory data analysis of individual data sources.* Before inferring the full model from all data, we explored the information content in each type of data individually (Figure 1-B,C,D).

*Protein-protein interactions (PPI).* We compared data from several complementary high-throughput assays, all available from BioGRID (Stark et al., 2006) as well as a literature-curated dataset (Reguly et al., 2006). We analyzed the overlap between the protein interactions and the pathway hypothesis of Fig 1-A. None of the datasets are informative for the membrane and nuclear compartments. Surprisingly, in the cytoplasm compartment we found that all of the high-throughput datasets show only  $\leq 3$  interactions between any of the proteins in the pathway. The situation was very different for the literature-curated data. Here, 45 interactions in the cytoplasm compartment covered 22 out of the 28 edges there (sensitivity  $> 78\%$ , specificity  $> 87\%$ , see Fig 1-B1).

*TF-DNA binding data.* We used the transcription factor binding data of (Harbison et al., 2004), which was not used to define the TF targets in the pathways hypothesis. However, the ROC in Figure 1-B4 shows this data contains a very clear signal that distinguishes the targets posited in the biological model from all other pathway genes.

*Co-expression data.* For gene expression data, we examined datasets in which the pathway genes showed a significant difference in correlation structure from all other yeast genes (using the SPELL algorithm of (Hibbs et al., 2007)) resulting in 20 datasets from 15 publications (including Roberts et al., 2000; Gasch et al., 2000; Brem and Kruglyak, 2005). Figure 1-B2 shows ROCs for predicting edges in the nuclear compartment for all datasets (grey lines) and the concatenated data (black line). No curve improves much on random prediction (the main diagonal). The reason is biological: Because expression data are a poor surrogate for protein activity, TFs are often less well correlated to their targets than the targets are between each other (Figure 1-B3). For STE12, which regulates itself, all correlation coefficients exhibit a strong trend towards high positive correlation. Whereas MCM1, which is not self-regulating, is far less strongly correlated to its targets than the targets are between each other. Thus, in general it is more informative to use the correlation between targets for inference, which is consistently high whether or not a TF is transcriptionally regulated itself.

*Gene perturbation data.* Paths in the graph are visible in cause-effect datasets (Hughes et al., 2000; Roberts et al., 2000). We find only very small effects of perturbations in the pathway on the expression of members of the membrane and cytoplasm compartment including TFs. Figure 1-C summarizes this result for the Roberts et al. (2000) data. Very similar results were found for the Hughes et al. (2000) data. The four boxes correspond to the three compartments plus TFs. In each box, a vertical line corresponds to a perturbation in the pathway (some replicated). The dots show the fold-changes of the pathway genes in this compartment. Only in the nuclear compartment are wide-spread large fold-changes visible. This observation motivates the construction of our likelihood around the presence of paths between the knockout and genes in the nuclear compartment (see section 2). In this way, when the knockout is far enough upstream, there is information about edges in the cytoplasm as well, even if the proteins there show no effect on the transcriptional level.

*Protein phosphorylation.* A first example of node information is protein phosphorylation. The study of Gruhler et al. (2005) assessed differential phosphorylation of proteins in response to pheromone. Figure 1-D1 shows the log-ratios between the pheromone treated and untreated conditions. Almost all proteins of the pheromone pathway measured by Gruhler et al. (2005) are up-regulated, which makes sense for a kinase cascade. The phosphorylation we observe for proteins corresponding to genes only attributed to the nuclear compartment in our model must be due to other kinase pathways in the cell. We further assessed to what extent the differential phosphorylation is correlated with the pathway model by fitting an auto-logistic regression. As a measure of correlation we computed the variance explained,  $R^2 = 0.76$ , using the bootstrap. The variance explained by the auto-logistic regression was found statistically significant, when compared to the correlation of differential phosphorylation with randomized pathway models,  $p \approx 0.062$ , and with randomized protein permutations on the true pathway model,  $p \approx 0.059$ .

*Protein domains.* A second example of node information are protein domains. We retrieved protein domains from PFAM (Punta et al., 2012) and SMART (Letunic et al., 2012). First, we sought to quantify which domains, if any, were over-represented in the set of proteins involved in the complete pheromone response pathway as well as in each compartment, in turn. Figure 1-D2 lists the domains that were found to be over-represented in the complete pathway and in the cytoplasm; darker shades of gray indicate a more significant p-value for the over-representation test.

Second, we sought to quantify to what extent the presence or absence of specific protein domains in proteins interacting with a given protein,  $P$ , was informative about the presence or absence of the same domain in such protein,  $P$ . This analysis was carried out using auto-logistic models, which summarize the informativeness of protein domains between interacting proteins on average, across all proteins in a given pathway. We fit auto-logistic regressions using each protein  $P$  in the cytoplasm compartment of the pheromone response pathway as data point, and the presence or absence of domains  $D_{1:K}$  in any one protein among those interacting with  $P$  as covariates.

We fit multivariate models, which assume that the presence or absence of either the same or complementary domains is a factor that facilitates protein physical interactions. The two tables in 1-D3 summarize the goodness of fit of the multivariate models, and report bootstrap p-values to assess the significance of the AIC scores. Figure 1-D3 shows the p-values obtained by fitting the multivariate auto-logistic regression to randomized pathway models. The domains identified by the multivariate models as putatively carrying signal about the pheromone pathway in the cytoplasm overlap with the domains identified by the over-representation analysis above; namely, P21 rho-binding domains, S-TKc domains, and tyrosine-specific catalytic domains.

In summary, node attributes of the proteins involved in the pheromone response pathways are informative about mechanistic elements of the kinase cascade, across cellular localizations and in the cytoplasm. These findings suggest that integrating node attributes such as protein domains and cellular localization should increase the likelihood of pathway models that encode real biological signal about the inner working of a target pathway.

*Data Integration.* The previous results suggest that some datasets are indeed more informative in certain cellular locations. For example, protein interactions can explain wide parts of the kinase cascade in the cytoplasm, while co-expression is very strong for TF targets. However, no dataset is informative in all compartments: Neither protein interactions nor knockout data can explain a complete pathway. The pheromone response pathway is an archetypical MAPK pathway, so we expect these observations also to be valid for other MAPK and signaling pathways. These results suggest that the compartment-specific modeling approach we take here is sensible.

As a proof of concept, we use the results of exploratory data analysis to construct the compartment map,  $\mathcal{C}$  (Table 2). That is, we fix the compartment map based on basic biological principles and the above exploratory analysis (see Figure 1). We briefly explore a sensitivity analysis on the com-

partment map in the supplementary results (Figure 6).

3.2. *Validation of the integrative pathway refinement strategy.* We evaluated how well the joint model, which combines all the complementary data types discussed above, supports the pathway hypothesis in Section 3 by sampling 1000 possible pathways using MCMC and tabulating the posterior probabilities over the edges. We demonstrate reasonable MCMC convergence in Figure 5.

Note that the logistic regression model for domain data may be subject to over-fitting and separation. This can occur since there are many different protein domains present, yet the frequency of any single domain is fairly low. To mitigate this issue, we used a Cauchy prior on the coefficients for the auto-logistic regression, which is a sensible default prior for this model (Gelman, 2008). Since the domain information in the pheromone pathway is relatively sparse, we also collected protein domain data from other MAPK pathways and used the hypothesized structure of those pathways to help learn the regression coefficients. Figure 1A includes the posterior probabilities for the edges in our initial hypothesis.

Further, we used a *leave-one-out* strategy to evaluate the predictive power of our model. We evaluated 37 separate fits where each node was in turn left out of the training pathway. The edges connected to this node were propagated to the neighboring nodes of the left-out node. We left out the nodes rather than edges, because specifically leaving out edges is equivalent to assuming that we know there is no edge present. We needed to construct our model in a way that encodes ignorance about the presence of an edge. Leaving out the nodes, instead of the edges, is one way of being agnostic about the presence of edges attached to that node. Only the coefficients in the auto-logistic regression were learned from the pathway hypothesis, so only the node likelihoods were affected. Table 3 shows the posterior probabilities for edges (under simulations in which a node was removed from

TABLE 2

*The compartment map,  $\mathcal{C}$ , associates pathway compartments with those data types that are informative for such compartments. Prior information is informative for all compartments.*

	Membrane	Cytoplasm	Nucleus
PPI	1	1	0
TF	0	0	1
Expr	0	0	1
Kout	0	1	1
Node	0	1	0
Prior	1	1	1

the prior hypothesis pathway). This table presents posterior probabilities for edges involved in knockout experiments.

For comparison, we also fit the model to *in silico* data. We constructed the “true pathway” to match the hypothesized MAPK pheromone pathway of Figure 1A. That is, we fixed a pathway with the matching nodes and edges. We then generated *in silico* datasets from the models specified in Section 2. The one exception is the data generation for the node data.

Here, we generate the presence of domains in a way such that short chains in the pathway are more likely to share domains than are random non-neighboring nodes. Specifically, we randomly chose chains of length 1 to 4 and added a common “domain” to every node in that chain. In this way, the domain data realistically reflect the notion that genes sharing common protein domains are more likely to interact.

The *in silico* leave-one-out results are also given in Table 3 beside the results for the true data. Figure 2 shows the precision-recall curve averaged over 30 simulated datasets. As in the true data analysis, the results demonstrate high precision and recall, especially in the “nucleus” and “cytoplasm”. The “membrane” shows the worst precision-recall because we have the fewest informative data types there, but when simulating from the true

TABLE 3

*Posterior edge probabilities for leave-one-out trials involving edges in knockout experiments. Since we use a leave-node-out scheme, there are two posterior probabilities for an edge (corresponding to which of the two node endpoints were left out for that particular simulation).*

	Real data			In Silico		
	Min	Average	Max	Min	Average	Max
STE11/STE7	0.01	0.01	0.01	0.26	.31	0.36
MCM1/STE2	0.00	0.01	0.02	0.03	0.12	0.2
MF(ALPHA)1/STE2	0.00	0.00	0.01	0.01	0.19	0.36
FUS1/STE12	0.80	0.83	0.87	0.39	0.66	0.92
CDC42/STE18	0.00	0.00	0.00	0.00	0.16	0.31
FUS3/STE12	0.01	0.01	0.01	0.01	0.10	0.19
STE5/STE7	0.13	0.13	0.13	0.00	0.14	0.27
BN1/CDC42	0.49	0.55	0.61	0.20	0.24	0.28
FAR1/MCM1	0.00	0.00	0.00	0.24	0.26	0.27
FAR1/STE12	0.00	0.00	0.00	0.00	0.37	0.73
STE12/CHS1	0.80	0.82	0.83	0.01	0.02	0.03
STE12/FIG2	0.84	0.84	0.85	0.04	0.24	0.43
MCM1/AGA1	0.10	0.23	0.37	0.07	0.17	0.27
STE12/FIG1	0.00	0.00	0.00	0.42	0.70	0.98
STE12/CIK1	0.83	0.84	0.85	0.94	0.96	0.98
STE12/KAR5	0.83	0.83	0.84	0.23	0.30	0.37
STE12/GIC2	0.83	0.83	0.84	0.12	0.54	0.95
MCM1/SWI4	0.00	0.00	0.00	0.16	0.29	0.41



data generating process, we still do quite well.

Finally, Figure 2 shows the precision-recall curve for our model, by compartment. For the membrane compartment, only the PPI data is informative, and weakly so. Thus, it performs the most poorly, although there are also by far the fewest genes in this compartment. By contrast, the nuclear and cytoplasm compartments both have high precision and recall.

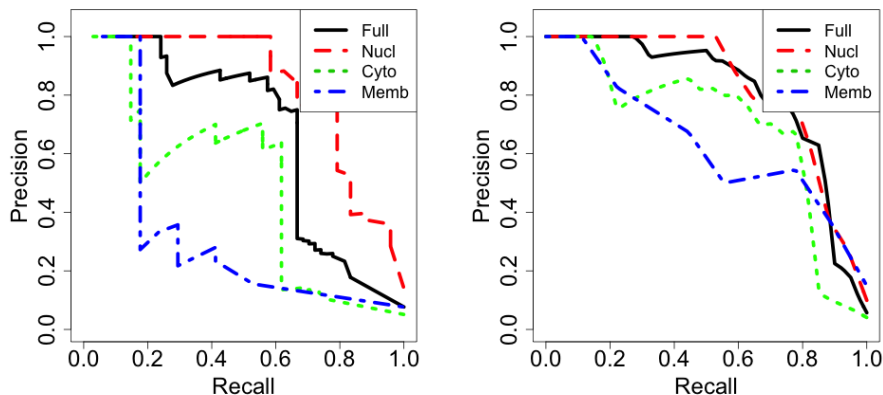


Fig 2: Precision/Recall curves overall and by compartment for the MAPK pathway (left) and simulated data (right). Thresholds are set on the posterior mean probability of an edge. In truth, the membrane compartment, which has the fewest genes, performs poorly because only the PPI dataset is (weakly) informative there. The simulated data reflects the average Precision/Recall over 30 simulated datasets (Section ).

*3.3. Inferring cross-talk with other pathways.* With our model, we are also able to identify possible cross-talk between pathways. In this paper, we focus on the pheromone response pathway, but our model can easily be used on other pathways, as long as we specify the relevant genes and transcription factors, and their corresponding cellular locations.

For instance, the MAPK pathway consists of the pheromone sub-pathway, as well as hypotonic shock, osmolarity and starvation sub-pathways. The degree of interaction between components of these MAPK pathways is not currently known. To identify cross-talk between the pheromone pathway and other MAPK pathways, we can simply include a new set of genes from the other sub pathways and fit the model as usual. The results for the cross-talk evaluations are displayed in Table 4.

**4. Discussion.** The proposed methodology achieves fairly strong predictive power by integrating data in a compartment specific way. Impor-

TABLE 4

*Number of inferred edges between the pheromone pathway and one of the other three sub-pathways with posterior probabilities above 0.3.*

	osmolarity	hypotonic	starvation
cytoplasm-cytoplasm	16	25	11
cytoplasm-membrane	12	17	8
cytoplasm-nucleus	22	17	3
cytoplasm-tf	0	2	3
membrane-membrane	2	2	2
membrane-nucleus	19	13	3
membrane-tf	0	1	2
nucleus-nucleus	4	7	0
nucleus-tf	1	6	10
tf-tf	0	0	2

tantly, we are able to evaluate how each data type contributes to the overall likelihood of any edge. Since each data type independently contributes to the probability of an edge, we can compute the fraction of the overall likelihood difference (between an edge and no edge) that is due to a particular data type. In this way our framework provides information about which parts of a pathway hypothesis are not well supported by available data (see Figure 3).

In addition, our methodology can identify if a particular data type tends to disagree with the other data types for sets of edges. This could indicate whether or not a data type is at all useful for modeling edges in a particular cellular location. A sensitivity analysis on the compartment map (Supplement) shows that indeed precision/recall degrades when non-informative data types are used to infer edges in certain cellular locations. Thus, it may be possible to do inference on the compartment map from Table 2, rather than fix it a priori. Alternatively, we could put a probabilistic prior over the entries in the compartment map which reflect any subjective uncertainty about where types are informative. Finally we could use this information can be used to check the validity of the individual data models of Section 2.

There are some open statistical issues that could be addressed in future work. For instance, one of the major challenges of pathway modeling is that typically we only track evidence for known edges and rarely record evidence for a lack of interaction. This makes supervised pathway inference very difficult. In our framework, we have strong a priori evidence for the presence of some edges edges, but no a priori evidence about the presence of non-edges. For this reason, in future work it may be a good idea to treat all edges without prior evidence of an interaction as missing data rather than a “true

zero”. Better documentation of experimentally verified non-interactions between gene products would also be very useful in future analyses.

Another problem relates to the sparsity of the protein domain data. While there is evidence of signal here, there may be an over-fitting problem. With more domain data, or perhaps broader domain categories, we may be able

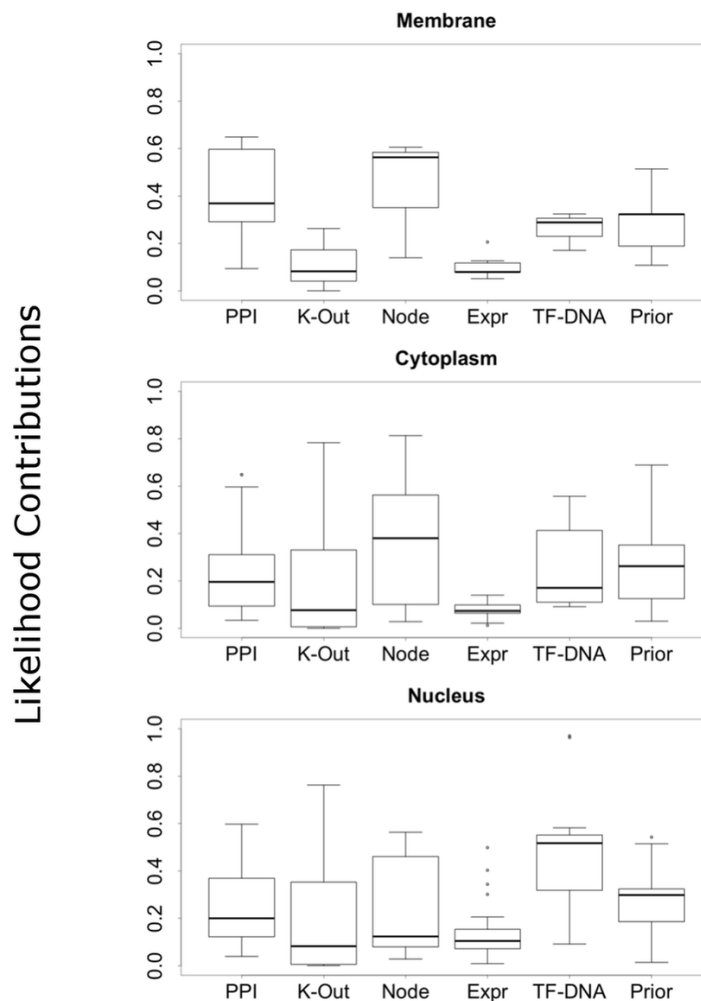


Fig 3: Percentages of differential likelihood (presence vs. absence of an edge) due to specific data types, by compartment. Node data contribute the most in the cytoplasm (center), whereas TF-DNA binding data contribute the most in the nucleus (bottom).

to learn more from the prior pathway. If this was the case, the leave-one-out results in the cytoplasm might improve significantly. This is evident from our results which show how borrowing domain information from other MAPK sub-pathways significantly improved the posterior probabilities of edges in the leave-one-out simulations.

We also noticed that most of the knockouts in the gene perturbation data set we used were generally downstream. If the knockouts were further upstream from perturbed genes in the nucleus, then we could learn about the possible presence of edges in a path between the knockout and other genes.

Lastly, we divided the pathway into its three main compartments: membrane, cytoplasm and nucleus. However, in future work, we hope to divide the pathway more finely into the over two dozen cellular components specified by the gene ontology (GO) for the yeast *S. Cerevisiae*. By dividing the pathway into more compartments, we would also have a greater degree of control over which data types are used in various parts of the cell.

4.1. *Concluding remarks.* In this paper we introduced a technique for refining cellular pathway models by integrating heterogeneous data sources in a compartment specific way and explicitly included node properties in our model. Our case-study results indicate that this model can be useful for discovering new components or cross-talk with other pathways. Our powerful and flexible pathway modeling framework can be easily extended and modified to include additional and novel datasets.

## APPENDIX A: SUPPLEMENTARY RESULTS

In this appendix we present convergence diagnostics for some network statistics, briefly explore sensitivity of the compartment map and present more details about the simulation results.

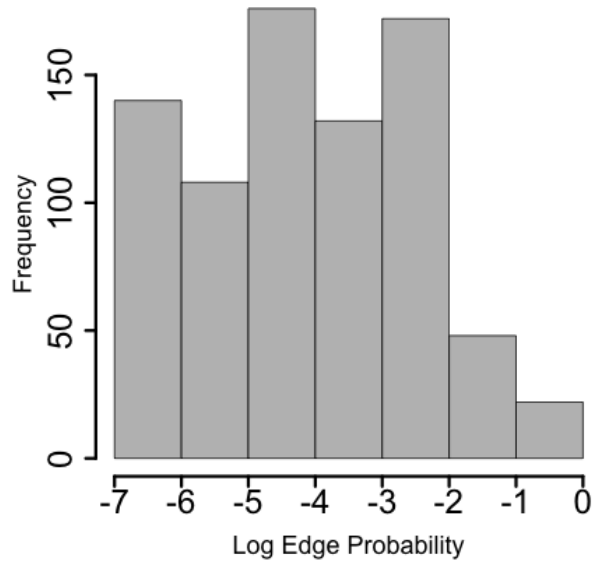


Fig 4: Log posterior probabilities for edges that were not in the hypothesis pathway. The vast majority of non-edges have small posterior probability (third quantile at 0.02). However, there are a few highly probable edges, which may indicate previously undiscovered interactions.

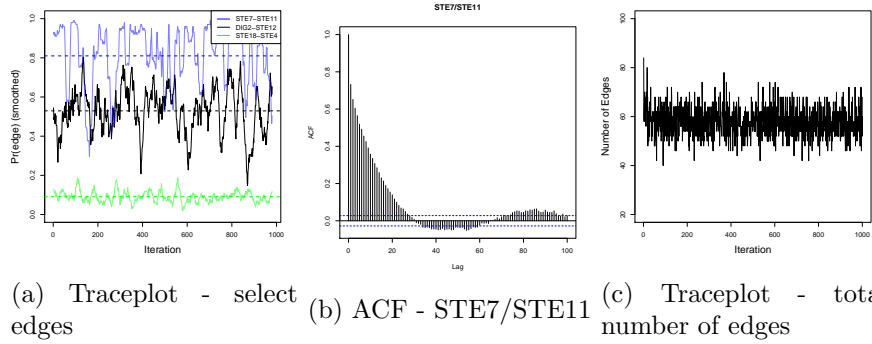


Fig 5: Convergence diagnostics from the fit to the MAPK pheromone response pathway. We initialize the sampler at the hypothesis pathway and find convergence to a stationary distribution at a local optima is achieved. a) Traceplot of kernel smoothed estimates of edge probabilities. b) Autocorrelation function for a select edge. c) Traceplot for total number of edges in the network. We calculate the Gelman-Rubin statistics for all parameters and find a potential scale reduction factor of less than 1.1 for all inferred edges. The effective sample size for the total number of edges in the network is about 500 (per 1000 saved samples) and the effective sample size for the presence of any individual edge in the network is closer to 100 (per 1000 samples).

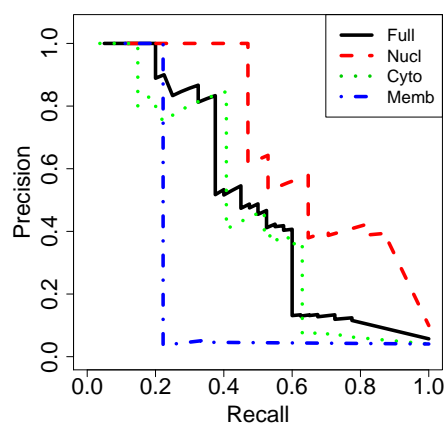


Fig 6: Precision/Recall curve for pathway inference using a prior on different compartment map initializations. The space of possible compartment initializations are constrained by basic biological principles and the models used in Section 3. However, within this space we posit a uniform prior on the compartment map. The results are comparable to those in Figure 1 although the precision and recall is generally slightly lower. This suggests that some map specifications actually add noise relative to the map chosen for the full analysis (Table 2)

	Gene 1	Gene 2	Prob		Gene 1	Gene 2	Prob
1	STE12	DIG2	0.60	61	FUS3	MSG5	0.05
2	STE12	FUS1	0.39	62	FUS3	FAR1	0.00
3	STE12	FUS3	0.01	63	MSG5	FUS3	0.00
4	STE12	FAR1	0.00	64	FAR1	STE12	0.73
5	STE12	MCM1	0.00	65	FAR1	FUS3	0.27
6	STE12	FIG2	0.43	66	FAR1	MCM1	0.27
7	STE12	FIG1	0.42	67	MCM1	STE12	0.00
8	STE12	CIK1	0.98	68	MCM1	MFA1	0.15
9	STE12	GIC2	0.12	69	MCM1	STE2	0.03
10	STE12	AFR1	0.01	70	MCM1	FAR1	0.24
11	STE12	KAR5	0.23	71	MCM1	SWI4	0.41
12	STE12	CHS1	0.03	72	MCM1	MFA2	0.20
13	STE12	AGA1	0.27	73	MCM1	AGA1	0.27
14	DIG2	STE12	0.68	74	MCM1	ALK1	0.15
15	DIG2	FUS3	0.00	75	MCM1	SWI5	0.38
16	STE7	STE11	0.26	76	MCM1	CDC20	0.34
17	STE7	STE5	0.21	77	SWI4	MCM1	0.16
18	STE7	FUS3	0.26	78	MFA2	MCM1	0.19
19	STE11	STE7	0.36	79	FIG2	STE12	0.04
20	STE11	STE20	0.24	80	FIG1	STE12	0.98
21	STE11	STE5	0.00	81	CIK1	STE12	0.94
22	STE20	STE11	0.00	82	GIC2	STE12	0.95
23	STE20	CDC42	0.31	83	AFR1	STE12	0.02
24	STE20	BEM1	0.08	84	KAR5	STE12	0.37
25	STE20	STE5	0.00	85	CHS1	STE12	0.01
26	CDC42	STE20	0.00	86	AGA1	STE12	0.00
27	CDC42	BNI1	0.28	87	AGA1	MCM1	0.07
28	CDC42	STE4	0.24	88	ALK1	MCM1	0.24
29	CDC42	STE18	0.31	89	SWI5	MCM1	0.13
30	CDC42	BEM1	0.47	90	CDC20	MCM1	0.18
31	CDC42	CDC24	0.50				
32	FUS1	STE12	0.98				
33	BNI1	CDC42	0.20				
34	MFA1	STE3	0.34				
35	MFA1	MCM1	0.07				
36	STE2	MF(ALPHA)2	0.01				
37	STE2	GPA1	0.35				
38	STE2	MCM1	0.20				
39	STE3	MFA1	0.30				
40	STE3	GPA1	0.13				
41	MF(ALPHA)2	STE2	0.36				
42	GPA1	STE2	0.01				
43	GPA1	STE3	0.14				
44	GPA1	STE4	0.14				
45	GPA1	STE18	0.12				
46	STE4	CDC42	0.22				
47	STE4	GPA1	0.14				
48	STE18	CDC42	0.00				
49	STE18	GPA1	0.13				
50	BEM1	STE20	0.36				
51	BEM1	CDC42	0.18				
52	CDC24	CDC42	0.18				
53	STE5	STE7	0.00				
54	STE5	STE11	0.00				
55	STE5	STE20	0.00				
56	STE5	FUS3	0.00				
57	FUS3	STE12	0.19				
58	FUS3	DIG2	0.21				
59	FUS3	STE7	0.22				
60	FUS3	STE5	0.05				

TABLE 5

*Posterior edge probabilities.*



## REFERENCES

- Alberts, B., A. Johnson, J. Lewis, M. Raff, K. Roberts, and P. Walter. (2002). *Molecular Biology of the Cell* (4 ed.). Garland Science. Available from: <http://www.ncbi.nlm.nih.gov/books/NBK21054/>.
- Balbin, O. A., J. R. Prensner, A. Sahu, A. Yocum, S. Shankar, R. Malik, D. Fermin, S. M. Dhanasekaran, B. Chandler, D. Thomas, D. G. Beer, X. Cao, A. I. Nesvizhskii, and A. M. Chinnaiyan (2013, Oct). Reconstructing targetable pathways in lung cancer by integrating diverse omics data. *Nat Commun* 4, 2617.
- Bernard, A. and A. J. Hartemink (2005). Informative structure priors: joint learning of dynamic regulatory networks from multiple types of data. *Pac Symp Biocomput*, 459–470.
- Brem, R. B. and L. Kruglyak (2005, Feb). The landscape of genetic complexity across 5,700 gene expression traits in yeast. *Proc Natl Acad Sci U S A* 102(5), 1572–1577.
- Friedman, N. (2004, Feb). Inferring cellular networks using probabilistic graphical models. *Science* 303(5659), 799–805.
- Fröhlich, H., T. Beißbarth, A. Tresch, D. Kostka, J. Jacob, R. Spang, and F. Markowetz (2008, Nov). Analyzing gene perturbation screens with nested effects models in r and bioconductor. *Bioinformatics* 24(21), 2549–2550.
- Fröhlich, H., M. Fellmann, H. Sülthmann, A. Poustka, and T. Beißbarth (2007). Large scale statistical inference of signaling pathways from rnaï and microarray data. *BMC Bioinformatics* 8, 386.
- Fröhlich, H., M. Fellmann, H. Sülthmann, A. Poustka, and T. Beißbarth (2008, Oct). Predicting pathway membership via domain signatures. *Bioinformatics* 24(19), 2137–2142.
- Gasch, A. P., P. T. Spellman, C. M. Kao, O. Carmel-Harel, M. B. Eisen, G. Storz, D. Botstein, and P. O. Brown (2000, Dec). Genomic expression programs in the response of yeast cells to environmental changes. *Mol Biol Cell* 11(12), 4241–4257.
- Gat-Viks, I. and R. Shamir (2007, Mar). Refinement and expansion of signaling pathways: the osmotic response network in yeast. *Genome Res* 17(3), 358–367.
- Gelman, A. (2008). A weakly informative default prior distribution for logistic and other regression models. *The Annals of Applied Statistics* 2(4), 1360–1383.
- Gitter, A., M. Carmi, N. Barkai, and Z. Bar-Joseph (2013, Feb). Linking the signaling cascades and dynamic regulatory networks controlling stress responses. *Genome Res* 23(2), 365–376.
- Gruhler, A., J. V. Olsen, S. Mohammed, P. Mortensen, N. J. Faergeman, M. Mann, and O. N. Jensen (2005, Mar). Quantitative phosphoproteomics applied to the yeast pheromone signaling pathway. *Mol Cell Proteomics* 4(3), 310–327.
- Guan, Y., D. Gorenshiteyn, M. Burmeister, A. K. Wong, J. C. Schimenti, M. A. Handel, C. J. Bult, M. A. Hibbs, and O. G. Troyanskaya (2012). Tissue-specific functional networks for prioritizing phenotype and disease genes. *PLoS Comput Biol* 8(9), e1002694.
- Guan, Y., C. L. Myers, D. C. Hess, Z. Barutcuoglu, A. A. Caudy, and O. G. Troyanskaya (2008). Predicting gene function in a hierarchical context with an ensemble of classifiers. *Genome Biol* 9 Suppl 1, S3.
- Hahne, F., A. Mehrle, D. Arlt, A. Poustka, S. Wiemann, and T. Beißbarth (2008). Extending pathways based on gene lists using InterPro domain signatures. *BMC Bioinformatics* 9, 3.
- Hara, K., T. Ono, K. Kuroda, and M. Ueda (2012, May). Membrane-displayed peptide ligand activates the pheromone response pathway in *saccharomyces cerevisiae*. *J Biochem* 151(5), 551–557.
- Harbison, C. T., D. B. Gordon, T. I. Lee, N. J. Rinaldi, K. D. Macisaac, T. W. Danford,

- N. M. Hannett, J.-B. Tagne, D. B. Reynolds, J. Yoo, E. G. Jennings, J. Zeitlinger, D. K. Pokholok, M. Kellis, P. A. Rolfe, K. T. Takusagawa, E. S. Lander, D. K. Gifford, E. Fraenkel, and R. A. Young (2004, Sep). Transcriptional regulatory code of a eukaryotic genome. *Nature* *431*(7004), 99–104.
- Hibbs, M. A., D. C. Hess, C. L. Myers, C. Huttenhower, K. Li, and O. G. Troyanskaya (2007, Oct). Exploring the functional landscape of gene expression: directed search of large microarray compendia. *Bioinformatics* *23*(20), 2692–2699.
- Hibbs, M. A., C. L. Myers, C. Huttenhower, D. C. Hess, K. Li, A. A. Caudy, and O. G. Troyanskaya (2008). Analysis of computational functional genomic approaches for directing experimental biology: a case study in mitochondrial inheritance. *PLoS Comput Biol* *in press*.
- Hughes, T. R., M. J. Marton, A. R. Jones, C. J. Roberts, R. Stoughton, C. D. Armour, H. A. Bennett, E. Coffey, H. Dai, Y. D. He, M. J. Kidd, A. M. King, M. R. Meyer, D. Slade, P. Y. Lum, S. B. Stepaniants, D. D. Shoemaker, D. Gachotte, K. Chakraborty, J. Simon, M. Bard, and S. H. Friend (2000, Jul). Functional discovery via a compendium of expression profiles. *Cell* *102*(1), 109–126.
- Hyduke, D. R. and B. . Palsson (2010, Apr). Towards genome-scale signalling network reconstructions. *Nat Rev Genet* *11*(4), 297–307.
- Isci, S., H. Dogan, C. Ozturk, and H. H. Otu (2013, Nov). Bayesian network prior: network analysis of biological data using external knowledge. *Bioinformatics*.
- Kanehisa, M. and S. Goto (2000, Jan). Kegg: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res* *28*(1), 27–30.
- Kirouac, D. C., J. Saez-Rodriguez, J. Swantek, J. M. Burke, D. A. Lauffenburger, and P. K. Sorger (2012). Creating and analyzing pathway and protein interaction compendia for modelling signal transduction networks. *BMC Syst Biol* *6*, 29.
- Knapp, B. and L. Kaderali (2013). Reconstruction of cellular signal transduction networks using perturbation assays and linear programming. *PLoS One* *8*(7), e69220.
- Kofahl, B. and E. Klipp (2004, Jul). Modelling the dynamics of the yeast pheromone pathway. *Yeast* *21*(10), 831–850.
- Letunic, I., T. Doerks, and P. Bork (2012, Jan). Smart 7: recent updates to the protein domain annotation resource. *Nucleic Acids Res* *40*(Database issue), D302–D305.
- Li, J., H. Wei, T. Liu, and P. X. Zhao (2013, Oct). Gplexus: enabling genome-scale gene association network reconstruction and analysis for very large-scale expression data. *Nucleic Acids Res*.
- Llewellyn, R. and D. S. Eisenberg (2008, Nov). Annotating proteins with generalized functional linkages. *Proc Natl Acad Sci U S A*.
- Lo, K., A. E. Raftery, K. M. Dombek, J. Zhu, E. E. Schadt, R. E. Bumgarner, and K. Y. Yeung (2012). Integrating external biological knowledge in the construction of regulatory networks from time-series expression data. *BMC Syst Biol* *6*, 101.
- Markowetz, F., D. Kostka, O. G. Troyanskaya, and R. Spang (2007, Jul). Nested effects models for high-dimensional phenotyping screens. *Bioinformatics* *23*(13), i305–i312.
- Markowetz, F. and R. Spang (2007). Inferring cellular networks—a review. *BMC Bioinformatics* *8 Suppl 6*, S5.
- McClellan, M. N., A. Mody, J. R. Broach, and S. Ramanathan (2007, Mar). Cross-talk and decision making in MAP kinase pathways. *Nat Genet* *39*(3), 409–414.
- Mukherjee, S. and T. P. Speed (2008, Sep). Network inference using informative priors. *Proc. Natl. Acad. Sci. U.S.A.* *105*(38), 14313–14318.
- Mulder, K. W., X. Wang, C. Escriu, Y. Ito, R. F. Schwarz, J. Gillis, G. Sirokmny, G. Donati, S. Uribe-Lewis, P. Pavlidis, A. Murrell, F. Markowetz, and F. M. Watt (2012, Jul). Diverse epigenetic strategies interact to control epidermal differentiation. *Nat*

- Cell Biol* 14(7), 753–763.
- Müller, P., D. Kутtenkeuler, V. Gesellchen, M. P. Zeidler, and M. Boutros (2005, Aug). Identification of JAK/STAT signalling components by genome-wide rna interference. *Nature* 436(7052), 871–875.
- Myers, C. L., D. Robson, A. Wible, M. A. Hibbs, C. Chiriac, C. L. Theesfeld, K. Dolinski, and O. G. Troyanskaya (2005). Discovery of biological networks from diverse functional genomic data. *Genome Biol* 6(13), R114.
- Nagiec, M. J. and H. G. Dohlman (2012, Jan). Checkpoints in a yeast differentiation pathway coordinate signaling during hyperosmotic stress. *PLoS Genet* 8(1), e1002437.
- Nariai, N., S. Kim, S. Imoto, and S. Miyano (2004). Using protein-protein interactions for refining gene networks estimated from microarray data by bayesian networks. *Pac Symp Biocomput*, 336–347.
- Ourfali, O., T. Shlomi, T. Ideker, E. Ruppın, and R. Sharan (2007, Jul). SPINE: a framework for signaling-regulatory pathway inference from cause-effect experiments. *Bioinformatics* 23(13), i359–i366.
- Pounds, S. and S. W. Morris (2003). Estimating the occurrence of false positives and false negatives in microarray studies by approximating and partitioning the empirical distribution of p-values. *Bioinformatics* 19(10), 1236–1242.
- Punta, M., P. C. Coghill, R. Y. Eberhardt, J. Mistry, J. Tate, C. Bournsell, N. Pang, K. Forslund, G. Ceric, J. Clements, A. Heger, L. Holm, E. L. L. Sonnhammer, S. R. Eddy, A. Bateman, and R. D. Finn (2012, Jan). The pfam protein families database. *Nucleic Acids Res* 40(Database issue), D290–D301.
- Reguly, T., A. Breitkreutz, L. Boucher, B.-J. Breitkreutz, G. C. Hon, C. L. Myers, A. Parsons, H. Friesen, R. Oughtred, A. Tong, C. Stark, Y. Ho, D. Botstein, B. Andrews, C. Boone, O. G. Troyanskaya, T. Ideker, K. Dolinski, N. N. Batada, and M. Tyers (2006). Comprehensive curation and analysis of global interaction networks in *saccharomyces cerevisiae*. *J Biol* 5(4), 11.
- Ren, B., F. Robert, J. J. Wyrick, O. Aparicio, E. G. Jennings, I. Simon, J. Zeitlinger, J. Schreiber, N. Hannett, E. Kanin, T. L. Volkert, C. J. Wilson, S. P. Bell, and R. A. Young (2000, Dec). Genome-wide location and function of dna binding proteins. *Science* 290(5500), 2306–2309.
- Roberts, C. J., B. Nelson, M. J. Marton, R. Stoughton, M. R. Meyer, H. A. Bennett, Y. D. He, H. Dai, W. L. Walker, T. R. Hughes, M. Tyers, C. Boone, and S. H. Friend (2000, Feb). Signaling and circuitry of multiple MAPK pathways revealed by a matrix of global gene expression profiles. *Science* 287(5454), 873–880.
- Ryan, C. J., P. Cimermani, Z. A. Szpiech, A. Sali, R. D. Hernandez, and N. J. Krogan (2013, Dec). High-resolution network biology: connecting sequence with function. *Nat Rev Genet* 14(12), 865–879.
- Schäfer, J. and K. Strimmer (2005a, Mar). An empirical Bayes approach to inferring large-scale gene association networks. *Bioinformatics* 21(6), 754–764.
- Schäfer, J. and K. Strimmer (2005b). A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Stat Appl Genet Mol Biol* 4, Article32.
- Schultz, J., F. Milpetz, P. Bork, and C. P. Ponting (1998, May). Smart, a simple modular architecture research tool: identification of signaling domains. *Proc Natl Acad Sci U S A* 95(11), 5857–5864.
- Scott, J., T. Ideker, R. M. Karp, and R. Sharan (2006, Mar). Efficient algorithms for detecting signaling pathways in protein interaction networks. *J Comput Biol* 13(2), 133–144.
- Segal, E., M. Shapira, A. Regev, D. Pe’er, D. Botstein, D. Koller, and N. Friedman (2003,

- Jun). Module networks: identifying regulatory modules and their condition-specific regulators from gene expression data. *Nat Genet* 34(2), 166–176.
- Segal, E., H. Wang, and D. Koller (2003). Discovering molecular pathways from protein interaction and gene expression data. *Bioinformatics* 19 Suppl 1, i264–i271.
- Simon, I., J. Barnett, N. Hannett, C. T. Harbison, N. J. Rinaldi, T. L. Volkert, J. J. Wyrick, J. Zeitlinger, D. K. Gifford, T. S. Jaakkola, and R. A. Young (2001, Sep). Serial regulation of transcriptional regulators in the yeast cell cycle. *Cell* 106(6), 697–708.
- Stark, C., B.-J. Breitkreutz, T. Reguly, L. Boucher, A. Breitkreutz, and M. Tyers (2006, Jan). BioGRID: a general repository for interaction datasets. *Nucleic Acids Res* 34(Database issue), D535–D539.
- Stelnic-Klotz, I., S. Legewie, O. Tchernitsa, F. Witzel, B. Klinger, C. Sers, H. Herzog, N. Blthgen, and R. Schfer (2012). Reverse engineering a hierarchical regulatory network downstream of oncogenic kras. *Mol Syst Biol* 8, 601.
- Tresch, A. and F. Markowetz (2008). Structure learning in nested effects models. *Stat Appl Genet Mol Biol* 7, Article9.
- Vaga, S., M. Bernardo-Faura, T. Cokelaer, A. Maiolica, C. A. Barnes, L. C. Gillet, B. Hege-  
mann, F. van Drogen, H. Sharifian, E. Klipp, M. Peter, J. Saez-Rodriguez, and R. Ae-  
bersold (2014). Phosphoproteomic analyses reveal novel cross-modulation mechanisms  
between two signaling pathways in yeast. *Mol Syst Biol* 10(12), 767.
- Wang, X., M. A. Castro, K. W. Mulder, and F. Markowetz (2012). Posterior association  
networks and functional modules inferred from rich phenotypes of gene perturbations.  
*PLoS Comput Biol* 8(6), e1002566.
- Wang, X., K. Yuan, C. Hellmayr, W. Liu, and F. Markowetz (2014). Reconstructing  
evolving signaling networks by hidden markov nested effects models. *Annals of Applied  
Statistics* 8(1), 1–647.
- Werhli, A. V. and D. Husmeier (2007). Reconstructing gene regulatory networks with  
Bayesian networks by combining expression data with multiple sources of prior knowl-  
edge. *Stat Appl Genet Mol Biol* 6, Article15.
- Workman, C. T., H. C. Mak, S. McCuine, J.-B. Tagne, M. Agarwal, O. Ozier, T. J. Begley,  
L. D. Samson, and T. Ideker (2006, May). A systems approach to mapping dna damage  
response pathways. *Science* 312(5776), 1054–1059.
- Yates, P. D. and N. D. Mukhopadhyay (2013). An inferential framework for biological  
network hypothesis tests. *BMC Bioinformatics* 14, 94.
- Yeang, C.-H., H. C. Mak, S. McCuine, C. Workman, T. Jaakkola, and T. Ideker (2005).  
Validation and refinement of gene-regulatory pathways on a network of physical inter-  
actions. *Genome Biol* 6(7), R62.
- Yip, K. Y., R. P. Alexander, K.-K. Yan, and M. Gerstein (2010). Improved reconstruction  
of in silico gene regulatory networks by integrating knockout and perturbation data.  
*PLoS One* 5(1), e8121.

§ HARVARD UNIVERSITY  
DEPARTMENT OF STATISTICS  
1 OXFORD STREET  
CAMBRIDGE, MA 02138, USA  
E-MAIL: [afranks@fas.harvard.edu](mailto:afranks@fas.harvard.edu)  
E-MAIL: [airoldi@fas.harvard.edu](mailto:airoldi@fas.harvard.edu)

¶ UNIVERSITY OF CAMBRIDGE  
CANCER RESEARCH UK CAMBRIDGE INSTITUTE  
LI KA SHING CENTRE  
ROBINSON WAY  
CAMBRIDGE, CB2 0RE, UK  
E-MAIL: [florian.markowetz@cruk.cam.ac.uk](mailto:florian.markowetz@cruk.cam.ac.uk)