

From correlation to causality: statistical approaches to learning regulatory relationships in large-scale biomolecular investigations

Robert O. Ness,^{*,†,‡} Karen Sachs,[¶] and Olga Vitek[‡]

Purdue University Department of Statistics, West Lafayette, College of Science, College of Computer and Information Science, Northeastern University, Boston, and School of Medicine, Stanford University, Palo Alto

E-mail: nessr@purdue.edu

Abstract

Causal inference – the task of uncovering regulatory relationships between components of biomolecular pathways and networks - is a primary goal of many high throughput investigations. Statistical associations between quantitative measurements can reveal an enticing number of putative causal interactions, but when do such associations reflect the underlying causal biomolecular mechanisms? The goal of this perspective is to provide suggestions for causal inference in large scale experiments, which utilize high throughput technologies such as mass spectrometry-based proteomics. We describe in non-technical terms the pitfalls of inference in large datasets, and suggest methods to overcome these pitfalls and reliably find regulatory associations.

^{*}To whom correspondence should be addressed

[†]Purdue University

[‡]Northeastern University

[¶]Stanford University

Introduction

Causal inference¹ applied to 'omics experiments, elucidates statistical associations [Added or edited: (e.g. Pearson correlation, Spearman correlation, mutual information)], which result from the underlying biomolecular mechanistic relationships. Modern high-throughput technologies such as mass spectrometry-based proteomics quantify components of the biomolecular systems on a large scale. [Added or edited: Thus, the opportunity exists for the proteomics community to apply computational methods for causal inference developed for gene microarray data to proteomics data.²]

We will demonstrate the paradoxical result that the large amount of data generated by high-throughput proteomics experiments make the task of causal inference more difficult. If only two biologically unrelated analytes were measured simultaneously, there is some probability of a spurious association – an association between the analytes arising purely from random chance – though the chance of this happening is low and gets lower as sample size increases. The difficulty of high-throughput data stems from the fact that the more analytes that are measured, the more pairwise associations between analytes are generated, resulting in a much greater chance for the appearance of spurious associations that are just as strong as associations that represent actual biological relationships. This obfuscates the true signal, leading to increased false discoveries of putative causal events. [Added or edited: While the gene network inference community has developed approaches to dealing with this issue², the solutions specific to the experimental conditionals and biological questions relevant to proteomics are still in early stages.] The problem can be addressed by refining the biological question, and by improving experimental design in terms of selection of (1) the subset of analytes, (2) the number of biological replicates, and (3) the type of biological conditions and stresses. Below, we describe in non-technical terms the process of elucidating causal associations from high-throughput data and suggest practical approaches for analysis of large scale [Added or edited: proteomics] datasets.

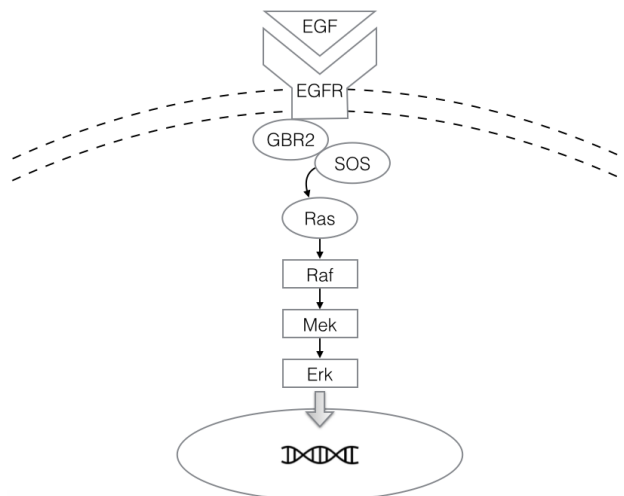


Figure 1: EGFR MAPK signaling pathway, an example of a pathway containing the phosphorylation cascade from Raf to Mek to Erk. The binding of ligand EGF to EGFR initiates a signal that leads to the cascade, which in turn regulates transcription. This cascade implies two direct causal relationships, namely $\text{Raf} \rightarrow \text{MEK}$, and $\text{Mek} \rightarrow \text{Erk}$. Raf and Erk have an indirect causal relationship, through Mek.

I. Small-scale statistical inference of causal relationships: conditional independence and interventions

Consider, e.g. the MAPK signaling cascade in Figure 1, which is part of several signaling pathways such as the EGFR MAPK pathway.³ In this cascade Raf causally affects the level of active (i.e., phosphorylated) Mek, while Mek causally affects Erk. Imagine these causal relationships were unknown: could they be detected from quantitative measurements on these phosphoproteins?

To illustrate the process of causal inference in this context, we simulated artificial data using the computational Huang-Ferrell model⁴ of this cascade. The model represents the key binding, phosphorylation, and dephosphorylation reactions of the cascade with mass action kinetics, and replicates the MAPK key signaling behavior observed in nature. We used the model to simulate an experiment with 50 replicate biological samples, and measurements of concentration (umol) of phosphorylated Raf, and doubly phosphorylated Mek and Erk in each sample.

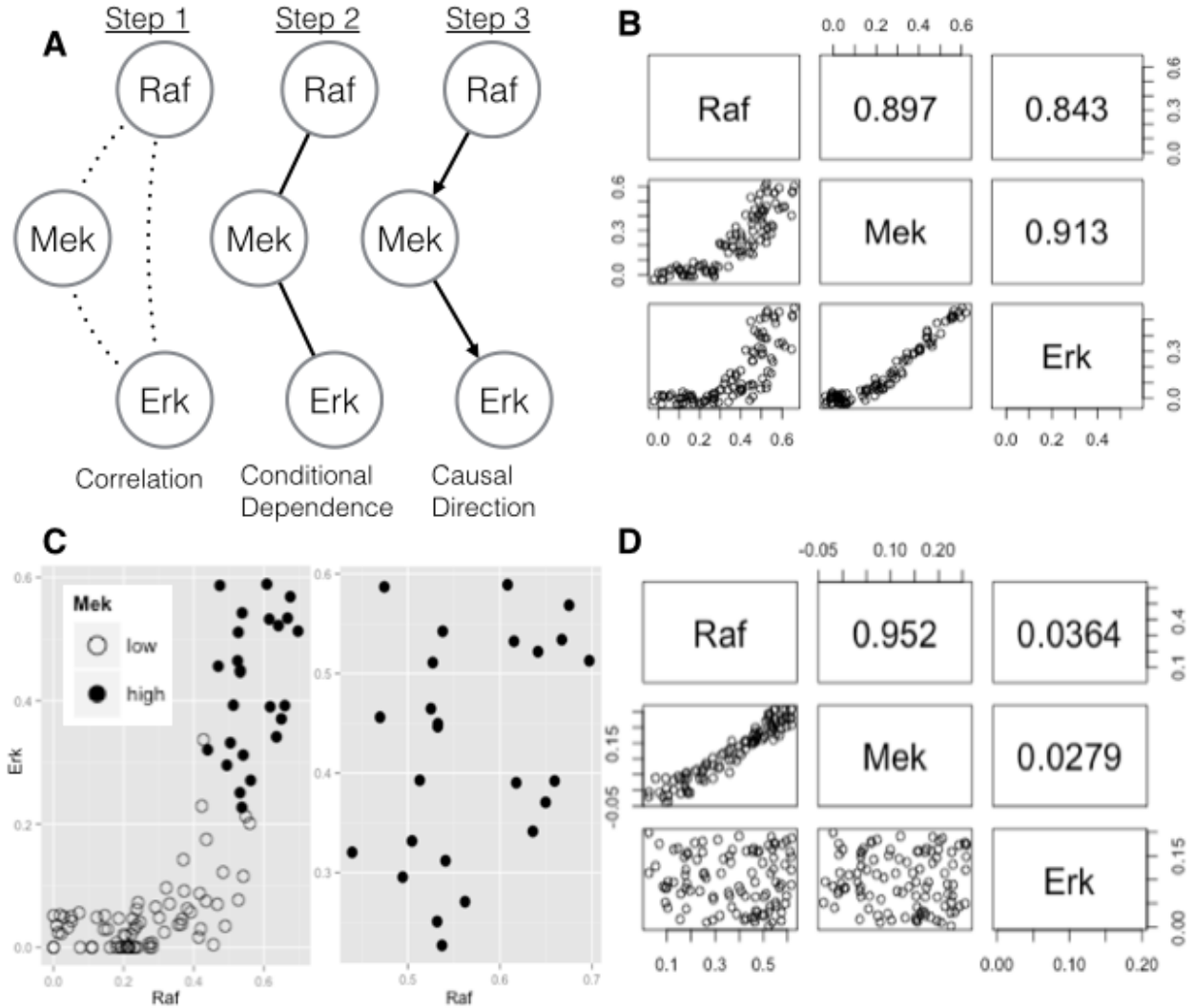


Figure 2: A: Overview of the 3 steps of causal inference, illustrated for the MAPK signaling cascade. B, C, and D feature an experiment simulated from the Huang-Ferrell computational model of the phosphorylation cascade. B: Pairwise plots of concentration values of phosphorylated (doubly phosphorylated for Mek and Erk) forms of each protein, and observed Spearman correlation values. The Raf – Erk correlation is high, despite the fact that Raf does not directly regulate Erk. C: Concentrations of Raf versus Erk, where samples corresponding to high Mek are highlighted in circles. The right panel of B shows the subset of samples with high Mek (i.e. conditional on Mek being high). In these samples the association between Raf and Erk disappears, indicating that Raf and Erk are *conditionally independent* given Mek. D: After Raf is inhibited, the observed association between Raf and Mek remains, while the association between Mek and Erk disappears. A second simulation, which is similar to the previous one but incorporates the inhibition of Mek reveals the causal flow from Raf to Mek to Erk.

Figure 2A demonstrates the causal inference workflow starting with analysis of statistical associations in the data. In step 1, a correlation graph between cascade components Raf, Mek, and Erk is assembled from the measurements of protein concentration. Step 2 reduces the correlation graph to a sparse graph of conditional dependencies (Raf–Mek, and Mek–Erk). Step 3 interrogates this graph to find putative causal relationships (Raf→Mek, and Mek→Erk). While step 1 has little requirements, step 2 requires multiple samples or replicates, and step 3 requires systematic interventions (e.g. with protein inhibitors).

Figure 2B illustrates Step 1 of the causal inference, and shows 2-way plots of the protein concentrations across the biological samples, and Spearman correlations to quantify the extent of the associations. The correlation values are high, and would meet most reasonable cut-off thresholds for constructing the correlation network in the left part of panel A. The Raf–Mek and the Mek–Erk correlation edges match the Raf→Mek, Mek→Erk known causal edges. What about the noncausal Raf–Erk edge? Despite the high Raf–Erk correlation, there is no direct causal mechanism between them (aside from the one via Mek, which is already accounted for via the Raf→Mek and Mek→Erk edges). In causal inference, our goal is to eliminate this "nuisance" edge. How is this done?

To describe Step 2 of causal inference, we introduce some terminology. In language of [Added or edited: probability] the quantified proteins are called *variables*. [Added or edited: When two proteins are connected in how they vary, such that knowing the behavior of one protein provides information on the other, they are called *dependent*. If there is no such connection, they are called *independent*. *Conditional independence* is a special case that is important to causal inference. Two dependent variables are conditionally independent if, after knowing, or *conditioning on* the behavior of third-party variables, they become independent – knowing one provides no additional information about the other after the information from the third party variables is considered.]

[Added or edited: Causation and dependence are linked as follows;

- For two variables, if one is the direct cause of the other, they are dependent, even after

conditioning on other variables in the system.

- If the two variables that are not directly causal related, *more often than not are still dependent*, but that dependence disappears after conditioning on the other variables in the system.

Note *dependence* is a property of the underlying biology, and *statistical association* is a property of the data generated by the biology, they are not the same. *Under the right measurement conditions* dependence manifests itself as statistical association in the data and application of statistical analysis can reveal cases of conditional independence (when those conditions breakdown is beyond the scope of this perspective, but Sachs and Itani et al. 2013⁵ explores this in detail within the phospho-proteomics context). So if we had an oracle that revealed that two proteins are still dependent after conditioning on all others, then we'd know they were one of them was a direct causal of the other. Without an oracle, we must rely on statistical inference of cases where pairs of proteins were conditional independent, ignore those cases and focus on the remaining pairs as potential causal relations. For a more rigorous treatment of the link between mechanism, probabilistic dependence, and statistics, consult Pearl's seminal work in this area.^{1]}

Let's see how this applies to the MAPK signaling cascade. [Added or edited: We know from the biology that Raf regulates the abundance of Erk through Mek. Thus Raf and Erk are dependent. However Erk abundance is directly determined by Mek, Raf does not exert any additional regulatory influence on Erk beyond its influence *through Mek*. Thus Raf and Erk are conditionally independent given Mek. This biology is manifests itself in the data.] Figure 2C compares Raf to Erk, indicating in [Added or edited: solid circles] the biological samples where concentrations of Mek are high (here, set to the top quartile). Note that when we subset the data to only the samples with high Mek, we can no longer visually detect the association between Raf and Erk. [Added or edited: To understand why we'd interpret this as evidence for the presence of conditional independence between Raf and Mek, we can build intuition from a formal statistical test of conditional independence – the null hypothesis is

conditional independence between Raf and Erk given Mek, the quantification of association between Raf and Erk *after subsetting Mek* is the test statistic, and we decide in favor of the null hypothesis if the p-value for the test statistic (the probability an association at least as high as the one observed could have emerged spuriously) is too high. Indeed, the class of causal inference algorithms called constraint-based algorithms repeatedly performs this type of test over multiple pairs of variables.⁶ Given this evidence the edge between Raf and Erk is removed, resulting in the middle graph of Figure 1A.]

At Step 2 in the figure, the direction of the regulation remains unknown. Inference of the direction of the chain of events requires that the experimental design involves external interventions or stresses. Figure 2D illustrates the results of Step 3, in the case where an intervention targeted Mek with an inhibitor. The intervention does not affect the concentration of Mek, however it blocks its ability to phosphorylate other proteins. After this intervention the Raf–Mek relationship is unchanged, while Erk drops to a low level. From this we can infer that Mek has causal influence on Erk, and since Raf was unaffected by the intervention, that Raf has causal influence on Mek. With the intervention, we can finally move from the undirected graph in panel A - step 2 to the causal graph in panel A - step 3.

In the general case, computational methods for causal inference follow the workflow in Figure 2A, while scaling it to characterize multiple inter-related variables. Step 1 creates a dense network of pairwise associations. [Added or edited: Step 2 identifies cases of conditional independence and removes those edges in those cases to create a much sparser network.] Finally, Step 3 uses the experimental design, specifically the information regarding the interventions, to evaluate these edges as evidence for potential causal events. See Koller-Friedman⁷ for a detailed description of these methods and their theoretical underpinnings. Numerous implementations of these algorithms are available in statistical software, e.g. in the R package `bnlearn`.⁸

II. Large-scale statistical inference of causal relationships: challenges of scaling up

A typical high-throughput experiment includes a small number of interventions, a small number of biological replicates, and quantifies a large number of analytes such as proteins. This creates challenges in each of the 3 steps of causal inference above.

In Step 1, the challenge is in quantifying statistical associations between each pair of the analytes across the biological samples. A large number of analytes yields a large number of spurious statistical associations, which arise without any biological justification, and are purely an artifact of random chance. Systematic pairwise relationships such as between Raf, Mek and Erk in the MAPK pathway will be obscured by the many spurious relationships that they will each form with causally unrelated proteins.

We illustrate this problem with a computer simulation, inspired by Fan et. al⁹ but translated to our context. First, we simulated an experiment that quantifies the abundances of 20 proteins in 100 biological samples. Second, we simulate another experiment where the number of proteins was increased to 500. In both experiments the proteins are completely independent from each other, and each protein in each replicate is assigned a value randomly drawn from a Gaussian distribution. In other words, we do not expect any biologically meaningful associations in these data. We repeated each of these simulations 500 times. Figure 3 shows for each experiment the histogram of the highest Pearson correlation across any pair of proteins in the 500 instances of the simulation. As can be seen, the experiment with 500 proteins produces relatively large maximum pairwise correlations, demonstrating that an increase in the number of proteins leads to an increase in spurious correlations. This is clearly a problem when high Pearson correlation is used as an initial evidence of a biological function.

Similarly, the increased incidence of spurious correlations impedes the performance of statistical methods in Step 2, which elucidate conditional independences in the data. If we

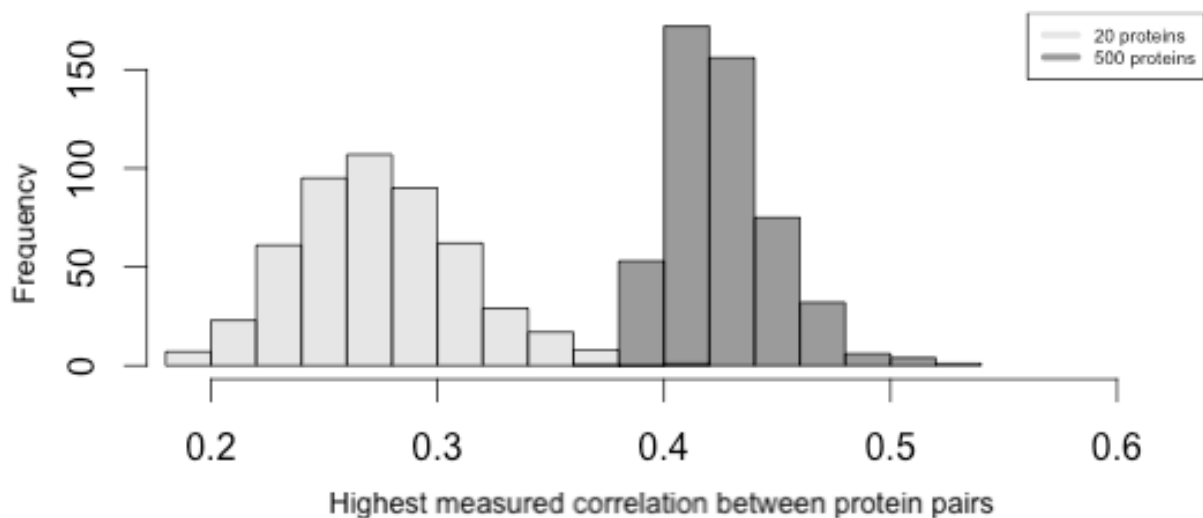


Figure 3: Simulated experiments quantifying 20 (or 500) unrelated proteins in 100 biological samples. The histograms display the highest Pearson correlation across any pair of proteins, calculated over 500 repetitions of the simulation. Increasing the number of proteins results in higher values of Pearson correlation without any biological justification.

think of this step as [Added or edited: attempting to detect which pairs are not conditionally independent], then the spurious correlations result in more false positives [Added or edited: detections]. To illustrate, we repeated the previous simulation, again starting with 20 proteins and 100 biological samples, but this time expanding to only 100 proteins. Instead of finding the highest spurious correlation between pairs of proteins, we apply an example of type of causal inference algorithm that performs a series of conditional independence tests between the sets of proteins (constraint-based algorithms⁶), and count the number of pairs of proteins [Added or edited: where the null hypothesis of conditional independence was rejected]. As before, since we randomly draw protein concentration measurement values from a Gaussian distribution, the values are completely independent, [Added or edited: no test should ever be rejected, and any instance of rejection by this algorithm would have biological justification]. We again repeat these experiments 500 times. Figure 4 shows the histograms of the counts of false positives. The figure graphs the results of the Grow-Shrink algorithm,¹⁰ [Added or edited: but we also repeated the simulation with the IAMB algorithm¹¹ and Inter-IAMB algorithm¹² and achieved similar results]. The results demonstrate that an increase

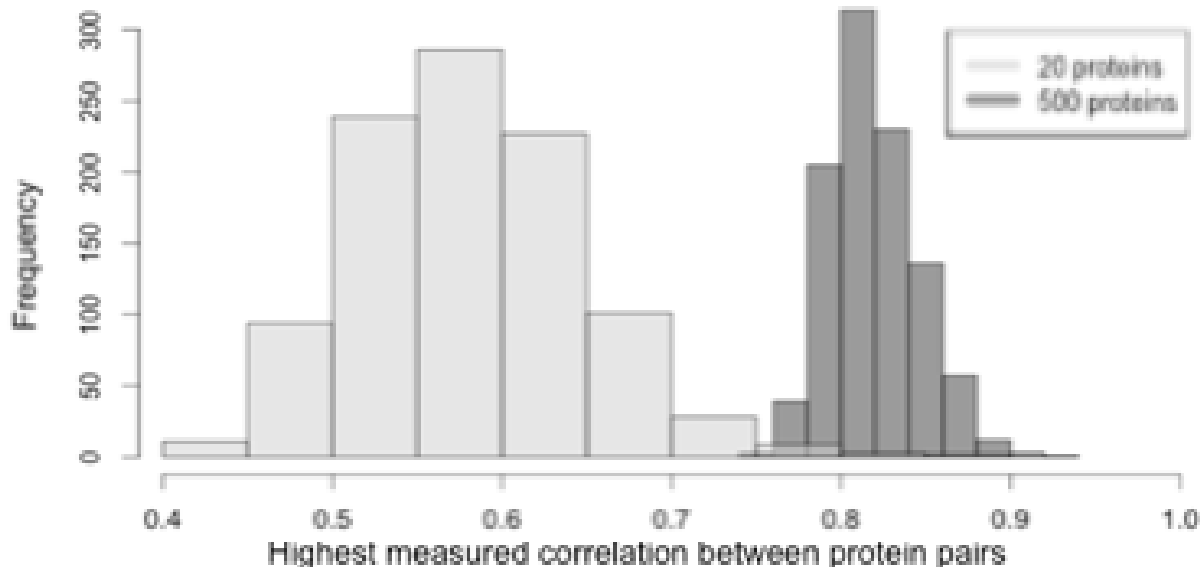


Figure 4: As in Figure 3, but reporting the counts of false positives (rejecting the null hypothesis of conditional independence).

in the number of proteins leads to an increase in false positives. This means that the computational methods for causal inference will fail for a typical large-scale experiment, because they cannot reliably distinguish the potential causal relationships from noise.

A related problem with high-throughput experiments in Step 2 is a relatively small number of biological replicates as compared to the number of analytes. While Step 1 (evaluating the set of pairwise associations) can be carried out with even a small amount of replicates, Step 2 relies on having many replicates relative to the number analytes. Note the experiments demonstrated in Figure 3 had 100 replicates to its 100 proteins, and still false positives were high. The quality of causal inference approaches decreases drastically as the numbers of replicates fall.

Finally, Step 3 of the workflow requires interventions to infer causal relationships from evidence of conditional independence. Quantifying a large number of analytes introduces a challenge for this step as well. As the number of features grows, the number of interventions needed to fully infer causality grows, and eventually performing a sufficient perturbation

experiments becomes infeasible. We showed with the Raf-Mek-Erk example that one intervention was sufficient to infer two causal relationships. [Added or edited: Given theoretical work on the number of interventions necessary to fully elucidate causality between a set of variations^{13 14 15}], a set of interventions on $< k$ variables may be sufficient to infer causal events between k analytes. Unfortunately, both the number of biological replicates and the number of interventions in high-throughput profiling experiments is typically small. Also, if the step 2 results in too many false positives for potential causal relations, this will adversely affect the results of step 3 regardless of the size of the dataset.

III. Approaches for inferring causality from high-throughput experiments

The problems outlined in section II paint a grim picture for causal inference in large datasets. Fortunately, these can be overcome, and effective causal inference can be a reality for large scale datasets. We provide suggestions for the best practices below.

1. *Limit the number of analytes.* Even though a list of analytes quantifiable with high-throughput technologies grows larger, only use a subset of measurements that are both biologically relevant and technologically accurate. The length of the list is not as important as the quality of measurements on the key parts of the system. If the broader biological system is well understood, it may be possible to design a targeted experiment that focuses on a specific network or pathway, and ask more specific questions of the data, such as the presence of a particular regulatory event. The more specific the question, the less data are needed to make solid causal conclusions.
2. *Profile more biological replicates.* The high-throughput measurements should provide more samples from distinct biological sources, which come from a same underlying population, in order to achieve a sufficient statistical power, and distinguish true and

spurious associations. This fact gives advantage to technologies, such as targeted proteomics tools, that quantify fewer analytes but have a higher sample throughput.

[Added or edited: A review by Twerve and Saez-Rodriguez provides a contrast of technologies by their content and sampling throughput.^{16]}

[Added or edited: Single cell mass cytometry is a technology deserving of special mention, where many thousands of cells per sample provide ample statistical power, though cell samples still need to be drawn from multiple individuals in order to make inferences about the population from which those individuals were drawn]. See Sachs et al 2005 for an in depth case study of network inference with a single cell dataset.¹⁷

3. *Use prior knowledge.* Prior knowledge improves the search for conditional independence and helps to determine causality. The prior knowledge can be in form of known canonical networks, extracted, e.g. from pathway databases such as KEGG. One example of such prior information is the MAPK pathway. The prior information reduces the search space of unknown associations that need to be considered, enables a more effective use of the data, and increases the confidence in newly discovered statistical associations. [Added or edited: Work by Saez-Rodriguez, Lauffenburger, and Sorger,¹⁸ as well as Terfve and Saez-Rodriguez¹⁹ use prior network knowledge to build logic models that reflect causal relationships between signaling proteins using proteomics measurements, the former using antibody/bead-based XMAP technology from Luminex, and the latter using LC-MS/MS.] Another example of prior knowledge is contextual information, such as spatial or temporal annotations of the quantitative measurements in the cell. The contextual information can be extracted from the literature or from other complementary (and potentially noisy) datasets. The causal inference algorithms can be extended to weigh evidence of conditional independence, depending on whether the analytes share the same spatial or temporal context.

4. *Select targeted interventions wisely.* Targeted interventions perturb individual compo-

nents of the biological system. An example is [Added or edited: siRNA knockdowns, as well as] small molecule inhibitors, which block the causal influence of a specific protein on its downstream components. Although effective, such targeted interventions are limited in number. Therefore, a strategic experimental design would use prior information, prioritize the interventions and the targets, and apply them to parts of the biological system that have most potential for new discovery of regulatory events. For example, a graph with undirected edges can be inspected, to reveal which nodes have potential to reveal the most causality if perturbed. Such targeted perturbations can be applied iteratively, after an initial statistical analysis revealed areas of the network where causal inference would benefit from extra measurements and data.

5. *Consider broad-scale interventions.* Broad-scale interventions sacrifice specificity of targeted interventions to simultaneously perturb many variables in a biological system. One example of broad-scale interventions is varying experimental conditions, in order to activate multiple pathways. Signals from endocrine, paracrine, and autocrine ligands elicit various signaling responses in hepatocytes, thus interventions that cover this range of signals gives the best picture of the broader causal network of hepatocyte signaling.²⁰ Similarly, interventions that go beyond receptor-level and perturb multiple components of the system bring cascading causal direct orientation deeper into the network. Although they do not provide specific information about the downstream effects of stimulation, broad-scale interventions can provide more causal insight. Therefore, the advantage of this approach is that it may enable elucidation of causality across the entire system.

This list suggests impactful approaches that can drastically improve causal inference from high-throughput experiments, by constraining the inference task, and thus allowing for accurate statistical inferences. For instance, the task of assessing which of all the possible KEGG pathways is present in a dataset will be far less error-prone than the task of assessing which of all possible combinations of my measured variables might form a biological pathway.

How should the tools listed be used? They are most powerful when used in combination, and in fact the lines between them are somewhat arbitrary and frequently blurred. For instance, using item #1 and item #2 in concert can be thought of as reducing the breadth and increasing the depth of the investigation. Items #4 and #5 call for use of interventions, but this task itself is complicated by measuring many analytes. Item #3, prior biological knowledge, can be used to prioritize what to target with that limited set of interventions. Causal inference becomes possible when using these tools in combination with a sound experimental design.

Acknowledgement

We acknowledge the participants of Dagstuhl seminar 15351 "Computational Mass Spectrometry" (December 2015) (<http://www.dagstuhl.de/de/programm/kalender/semhp/?semnr=15351>) for their contributions to the discussion on computational manuscripts.

References

- (1) Pearl, J. *Causality*; Cambridge university press, 2009.
- (2) Bensimon, A.; Heck, A. J.; Aebersold, R. *Annual review of biochemistry* **2012**, *81*, 379–405.
- (3) Holbro, T.; Hynes, N. E. *Annu. Rev. Pharmacol. Toxicol.* **2004**, *44*, 195–217.
- (4) Huang, C.-Y.; Ferrell, J. E. *Proceedings of the National Academy of Sciences* **1996**, *93*, 10078–10083.
- (5) Sachs, K.; Itani, S.; Fitzgerald, J.; Schoeberl, B.; Nolan, G.; Tomlin, C. *Interface focus* **2013**, *3*, 20130019.
- (6) Spirtes, P.; Glymour, C. N.; Scheines, R. *Causation, prediction, and search*; MIT press, 2000; Vol. 81.

- (7) Koller, D.; Friedman, N. *Probabilistic graphical models: principles and techniques*; MIT press, 2009.
- (8) Scutari, M. *arXiv preprint arXiv:0908.3817* **2009**,
- (9) Fan, J.; Han, F.; Liu, H. *National science review* **2014**, *1*, 293–314.
- (10) Margaritis, D. Learning Bayesian network model structure from data. Ph.D. thesis, US Army, 2003.
- (11) Tsamardinos, I.; Aliferis, C. F.; Statnikov, A. R.; Statnikov, E. Algorithms for Large Scale Markov Blanket Discovery. FLAIRS Conference. 2003.
- (12) Yaramakala, S.; Margaritis, D. Speculative Markov blanket discovery for optimal feature selection. Data mining, fifth IEEE international conference on. 2005; pp 4–pp.
- (13) Eberhardt, F.; Glymour, C.; Scheines, R. *arXiv preprint arXiv:1207.1389* **2012**,
- (14) Hauser, A.; Bühlmann, P. *arXiv preprint arXiv:1205.4174* **2012**,
- (15) Hyttinen, A.; Eberhardt, F.; Hoyer, P. O. *The Journal of Machine Learning Research* **2013**, *14*, 3041–3071.
- (16) Saez-Rodriguez, J.; Terfve, C. D. A. In *Modeling Signaling Networks Using High-throughput Phospho-proteomics* ; Goryanin, I. I., Goryachev, A. B., Eds.; Advances in Experimental Medicine and Biology; Springer New York: New York, NY, 2012; Vol. 736.
- (17) Sachs, K.; Perez, O.; Pe’er, D.; Lauffenburger, D. A.; Nolan, G. P. *Science* **2005**, *308*, 523–529.
- (18) Saez-Rodriguez, J.; Alexopoulos, L. G.; Epperlein, J.; Samaga, R.; Lauffenburger, D. A.; Klamt, S.; Sorger, P. K. *Molecular Systems Biology* **2009**, *5*, 331.

- (19) Terfve, C. D. A.; Wilkes, E. H.; Casado, P.; Cutillas, P. R.; Saez-Rodriguez, J. *Nature Communications* **2015**, *6*, 1–11.
- (20) Alexopoulos, L. G.; Saez-Rodriguez, J.; Cosgrove, B. D.; Lauffenburger, D. A.; Sorger, P. K. *Molecular & Cellular Proteomics* **2010**, *9*, 1849–1865.