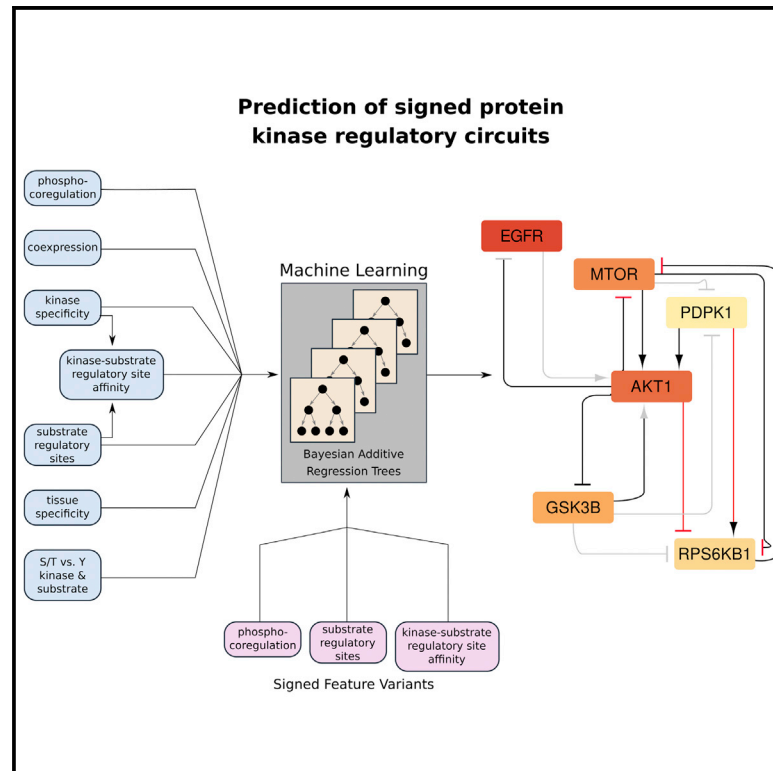


Prediction of Signed Protein Kinase Regulatory Circuits

Graphical Abstract



Highlights

- Combined diverse evidence for protein kinase-kinase regulatory relationships
- Used machine learning to predict (activating/inhibiting) regulatory relationships
- Recovered known signaling pathways and validated new relationships independently
- Analysis suggests that inter-kinase regulation is far denser than normally considered

Authors

Brandon M. Invergo,
Borghthor Petursson,
Nosheen Akhtar, ..., Pedro Cutillas,
Evangelia Petsalaki, Pedro Beltrao

Correspondence

p.cutillas@qmul.ac.uk (P.C.),
petsalaki@ebi.ac.uk (E.P.),
pbeltrao@ebi.ac.uk (P.B.)

In Brief

The activities of many protein kinases are themselves regulated by phosphorylation. While many human kinase-kinase regulatory relationships are known, the vast majority of potential relationships remains unexplored. Invergo et al. combined diverse data including phosphoproteomic, transcriptomic, and kinase specificity data in order to predict not only whether one human protein kinase regulates the activity of another, but also whether the relationship is activating or inhibiting. This model can reconstruct known signaling pathways, while also producing new predictions for experimental prioritization.



Article

Prediction of Signed Protein Kinase Regulatory Circuits

Brandon M. Invergo,^{1,3,4} Borgthor Petursson,^{1,3} Nosheen Akhtar,² David Bradley,¹ Girolamo Giudice,¹ Maruan Hijazi,² Pedro Cutillas,^{2,*} Evangelia Petsalaki,^{1,*} and Pedro Beltrao^{1,5,*}

¹European Molecular Biology Laboratory-European Bioinformatics Institute, Hinxton CB10 1SD, UK

²Centre for Genomics and Computational Biology, Queen Mary University of London, London EC1M 6BQ, UK

³These authors contributed equally

⁴Present address: Translational Research Exchange @ Exeter, University of Exeter, Exeter EX4 4QJ, UK

⁵Lead Contact

*Correspondence: p.cutillas@qmul.ac.uk (P.C.), petsalaki@ebi.ac.uk (E.P.), pbeltrao@ebi.ac.uk (P.B.)

<https://doi.org/10.1016/j.cels.2020.04.005>

SUMMARY

Complex networks of regulatory relationships between protein kinases comprise a major component of intracellular signaling. Although many kinase-kinase regulatory relationships have been described in detail, these tend to be limited to well-studied kinases whereas the majority of possible relationships remains unexplored. Here, we implement a data-driven, supervised machine learning method to predict human kinase-kinase regulatory relationships and whether they have activating or inhibiting effects. We incorporate high-throughput data, kinase specificity profiles, and structural information to produce our predictions. The results successfully recapitulate previously annotated regulatory relationships and can reconstruct known signaling pathways from the ground up. The full network of predictions is relatively sparse, with the vast majority of relationships assigned low probabilities. However, it nevertheless suggests denser modes of inter-kinase regulation than normally considered in intracellular signaling research. A record of this paper's transparent peer review process is included in the Supplemental Information.

INTRODUCTION

Cells continually respond and adapt to environmental stimuli. They employ sophisticated protein networks to propagate, amplify, and subsequently quench these signals efficiently. A common mechanism of relaying information from one protein to another is through reversible post-translational modifications (PTMs). Protein phosphorylation by kinases is one of the principal and best-studied PTMs. It plays a major role in cellular processes, such as growth, division, and differentiation (Acosta-Jaquez et al., 2009; Basson, 2012; Rhind and Russell, 2012).

Many protein kinases are themselves regulated by phosphorylation, giving rise to complex networks of kinase-kinase regulatory relationships. An accumulation of biochemical knowledge has produced consensus maps of several protein-kinase signaling pathways, which have been deposited in databases, such as Reactome (Fabregat et al., 2017), KEGG (Kanehisa et al., 2017), and SIGNOR (Perfetto et al., 2016). Kinase-kinase and other kinase-substrate relationships have also been annotated in databases, such as PhosphoSitePlus and Phospho.ELM (Dinkel et al., 2011; Hornbeck et al., 2015). However, a focus on well-studied protein kinases in the experimental investigation of kinase regulatory relationships overrepresents the activities of these kinases and has left the majority of the kinase-kinase interaction space largely unexplored (Invergo and Beltrao, 2018).

Similar effects have been reported for protein-protein interaction databases (Gillis et al., 2014). Subsequent proteome-wide analyses have found protein interactions to be ultimately more evenly spread across the proteome than previously indicated (Rolland et al., 2014), and the same is likely to be true for kinase signaling.

Incomplete maps of regulatory relationships could have serious impacts on systems-level analyses of signaling pathways. There is, therefore, a clear need for new methods for finding kinase-kinase regulatory relationships. Existing methods for data-driven reconstruction of signaling networks are generally designed for data that have been produced for the study of a specific pathway (e.g., via perturbation experiments) and typically benefit from the incorporation of prior knowledge about that pathway into the model (see, e.g., Hill et al., 2016; Invergo and Beltrao, 2018). The use of incomplete prior knowledge means that these methods are less likely to provide insight into broader patterns of protein-kinase regulation, especially of understudied kinases or cross-module signaling. However, recent advances in high-throughput phosphoproteomics, through liquid-chromatography tandem mass spectrometry (LC-MS/MS) and other technologies, show promise in the inference and analysis of signaling networks at larger scale (Babur et al., 2018; Rudolph et al., 2016; Terfve et al., 2015).



Alternatively, computational methods can be used to prioritize future experiments. Numerous previous attempts have been made to predict kinase-substrate relationships based on various features, such as amino acid sequences and/or functional information. Earlier methods, including Scansite (Obenauer et al., 2003) and NetPhosK (Blom et al., 2004), utilize position-specific scoring matrices (PSSMs) and neural networks to make predictions based on previously annotated substrates. Another method, GPS (group-based prediction system) (Zhou et al., 2004) bases prediction on peptide similarities and the Markov Cluster Algorithm. Later methods have integrated other features, such as probabilistic networks in addition to consensus sequences (Linding et al., 2007). Notably, however, these methods are geared toward the prediction of target phosphosites and do not make predictions about the regulatory impact of the phosphorylation. Some pathway reconstruction methods, such as modular response analysis (Kholodenko et al., 2002) or Bayesian techniques (Hill et al., 2012, 2017; Oates and Mukherjee, 2012; Oates et al., 2014), can infer kinase-kinase regulatory relationships, but they do not scale easily for many kinases, and they require purpose-built perturbation experimental data. Other methods scale to incorporate phosphoproteomic data for generalized predictions, but they require or benefit from the provision of a prior, literature-derived network from which to make predictions of regulatory relationships (see, e.g., Köksal et al., 2018; Rudolph et al., 2016; Terfve et al., 2015; Wilkes et al., 2015).

Here, we propose a supervised machine learning approach to estimate the probability of a functional, regulatory relationship between arbitrary pairs of human kinases, as well as to predict the sign (inhibiting or activating) of the regulation. We train the predictions on known kinase regulatory relationships by combining phosphoproteomic and transcriptomic data with kinase-substrate-sequence specificity models and a recently produced predictor of phosphosite functional impact (Ochoa et al., 2020). Our models allow us to make inferences even for kinases that lack any substrate annotations. The resulting network of predicted kinase-kinase regulatory relationships is highly modular and partitions into several clusters that reflect known functional associations, while suggesting denser modes of inter-regulation and feedback than typically assumed.

RESULTS

Regulatory Relationships Can Be Identified by Similar Phosphorylation Patterns at Functional Phosphosites and by Kinase Coexpression

We assume that kinases that are activated or inhibited in the same sets of conditions are more likely to be part of the same pathway and could form a regulatory interaction. Because many protein kinases are regulated by phosphorylation, we measured the correlation of phosphorylation of regulatory phosphosites for pairs of kinases across different conditions. If regulatory sites on two kinases show similar patterns of phosphorylation, one of the kinases might be responsible for regulating the other's activity. We assessed correlations of phosphosite quantification in two large-scale phosphoproteomic experiments (Mertins et al., 2016; Wilkes et al., 2015). Given that regulatory phosphosites have only been established for a small sub-

set of kinases, we employed a recently produced computational predictor of phosphosite functionality (Ochoa et al., 2020). This provided us with a score from 0.0 to 1.0 for each kinase phosphosite, with higher values indicating a stronger prediction of such sites regulating the kinase activity ("functional sites").

We found that kinase-kinase regulatory pairs often exhibit co-phosphorylation patterns at functional phosphosites. For example, mitogen-activated protein kinase 3 (MAPK3) is known to regulate the activity of ribosomal protein S6 kinases (Mérieu et al., 2000; Smith et al., 1999; Zhao et al., 1996). Indeed, we found strong correlation between functional sites T202 on MAPK3 and T577 on S6K-alpha-3 (RPS6KA3); meanwhile, no such correlation was found for atypical MAPK4, which has no known regulatory relationship with S6 kinases (Figure 1A). We quantified this relationship for each pair of sites between two kinases by producing a phosphosite "coregulation score," in which the log-transformed p value of the correlation is scaled by the two sites' functional scores (Figure 1A). We then checked whether known regulatory relationships annotated in the Omni-Path database (Türei et al., 2016) have higher coregulation scores than unannotated pairs. In both phosphoproteomic experiments, kinase-kinase regulatory pairs tend to exhibit higher maximum coregulation scores than pairs with no previously annotated relationship (one-sided Wilcoxon rank sum test, $W = 2.8 \times 10^7$, $p < 1 \times 10^{-6}$ [Mertins et al., 2016]; $W = 9.3 \times 10^5$, $p < 1 \times 10^{-6}$ [Wilkes et al., 2015]) (Figure 1B).

We next used two RNA sequencing (RNA-seq) datasets (GTEx Consortium, 2013; Uhlén et al., 2015) to test whether kinase co-expression is indicative of regulatory relationships. For example, if we consider the regulation of tyrosine-protein kinase BTK by Src-family protein kinases, we see a clear positive correlation between BTK expression and that of LYN (encoding tyrosine-protein kinase Lyn, a known regulator) (Cheng et al., 1994; Park et al., 1996; Rawlings et al., 1996). No such correlation exists for YES1 (tyrosine protein kinase Yes, which is not known to regulate BTK) (Figure 1C). In general, we found higher co-expression between pairs of kinases where a regulatory relationship exists than for those without any annotated relationship in both expression datasets (one-sided Wilcoxon rank sum test, $W = 1.3 \times 10^8$, $p < 1 \times 10^{-6}$ [GTEx Consortium, 2013]; $W = 1.3 \times 10^8$, $p < 1 \times 10^{-6}$ [Uhlén et al., 2015]) (Figure 1B).

We also found that tissue specificity, as represented by the skewness of expression values across tissue samples, is further indicative of kinase regulatory relationships. Continuing from the previous example, we can see that BTK and LYN both have skewed expression profiles (high expression in a few tissues), whereas YES1 has relatively even expression across tissues (Figure 1D). If we consider the absolute difference between tissue specificities for pairs of protein kinases, we find that pairs with regulatory relationships tend to have more similar expression profiles than those with no annotated relationship (one-sided Wilcoxon rank sum test, $W = 8.9 \times 10^7$, $p < 1 \times 10^{-6}$) (Figure 1B).

Linking Sequence Specificity to Phosphosite Functional Impact Identifies Direct Regulation of Protein-Kinase Activity

Kinases show preferences for phosphorylating some substrates over others, determined by the specific

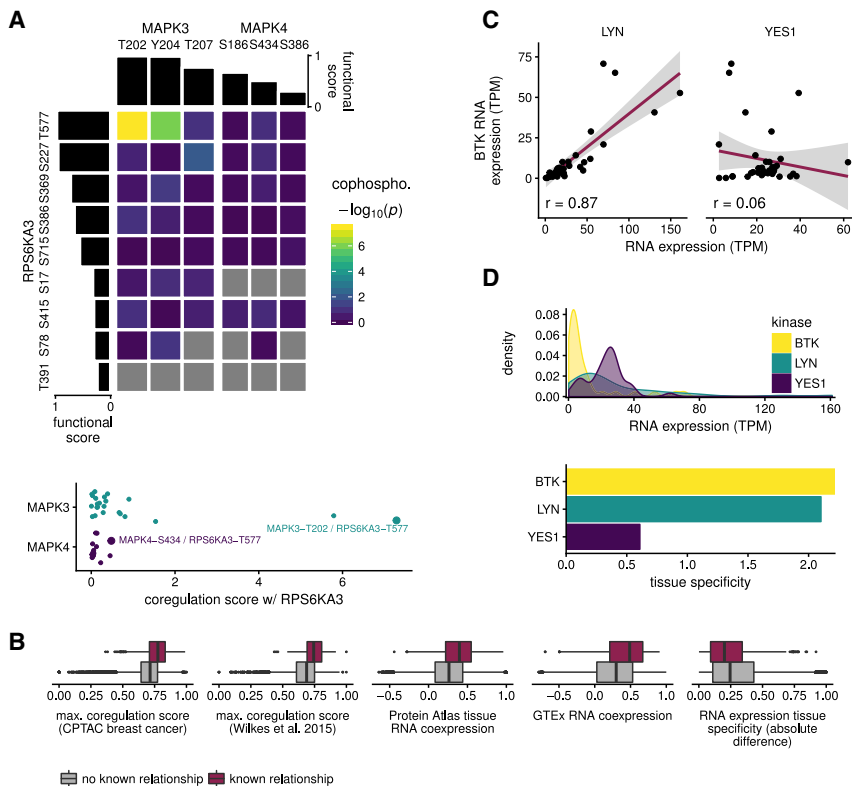


Figure 1. Correlations in Phosphorylation at Regulatory Sites or in Tissue Expression Patterns Are Predictive of Kinase-Kinase Regulatory Relationships

(A) Top: kinase MAPK3 exhibits cophosphorylation patterns at functional sites with RPS6KA3, a known substrate. The same patterns are not observed for MAPK4. Gray cells indicate missing values. Bottom: combining cophosphorylation p values and site functional scores provides an estimator of coregulation.

(B) Phospho-coregulation, tissue coexpression, and tissue specificity can discriminate cases of kinase-kinase regulation annotated in the OmniPath database from unannotated cases.

(C) The RNA transcripts encoding SRC-family kinase LYN and known substrate BTK show similar patterns of expression, while the expression of SRC-family kinase YES1, not known to regulate BTK, is unrelated. Shaded area represents 95% confidence intervals.

(D) Top: kernel-density estimates of the distributions of expression values across tissue samples for *BTK*, *LYN*, and *YES1*. Bottom: tissue specificity of RNA expression was quantified as the skewness of the kernel-density distributions. Here, *YES1* is more broadly expressed than the tissue-specific *LYN* and *BTK*.

phosphoacceptor residue and a surrounding amino acid sequence. By characterizing this specificity in a PSSM, we can score a kinase's potential for directly phosphorylating a putative substrate. However, we also wanted to determine, in an unbiased way, whether high-scoring substrate sites also tend to have regulatory effects. To achieve this, we employed the discounted cumulative gain (DCG) metric often used in the evaluation of information retrieval systems (Järvelin and Kekäläinen, 2002), wherein we treated a PSSM as a phosphosite "search function" and the functional score as a phosphosite "relevance metric."

Only 140 protein kinases had sufficient numbers of known substrate sites to build confident PSSMs. We have recently shown that proteins within the same kinase family tend to show similar specificity, which can be attributed to conserved specificity-determining residues (SDRs) within their protein-kinase domains (Bradley and Beltrao, 2019; Bradley et al., 2018). We thus investigated this as a means to assign PSSMs to kinases with insufficient substrate annotations. We first estimated the minimum residue similarity necessary across 10 kinase SDRs to make accurate PSSM assignments. We found that an SDR similarity of at least 0.8 (based on the BLOSUM62 amino acid substitution matrix) is needed to make assignments that are significantly better than a random assignment (Figure 2A). Nevertheless, this method of assignment did not substantially improve upon simply assigning a family-wise, composite PSSM (Figure 2B). Based on these results, we increased the coverage of kinases with PSSMs by assigning to under-annotated kinases a family-wise PSSM, where available ($n = 208$), or otherwise one via SDR similarity ($n = 14$), bringing the total

number of protein kinases with specificity profiles to 362 (Figure 2C).

Linking PSSM predictions to phosphosite functional scores via the DCG is best illustrated by an example. RAC-alpha serine/threonine-protein kinase (AKT1) has several phosphosites, a few of which have high functional scores. We consider two potential regulators: 3-phosphoinositide-dependent protein kinase 1 (PDPK1), a known regulator, and protein kinase C gamma type (PRKGC), not known to regulate AKT1. Some of AKT1's sites with the highest functional scores also score highly with PDPK1's PSSM, whereas PRKGC's PSSM favors sites with low functional scores (Figure 2D). These relationships can be quantified and visualized via the DCG: substrate sites are ranked by PSSM score and a cumulative sum of their functional scores is calculated, wherein each successive site contributes a smaller fraction of its functional score (Figure 2E). We can see that, although the two protein kinases achieve similar maximum PSSM scores, only PDPK1 produces a high DCG (Figure 2F).

As would be expected, we found that the PSSMs of known regulators in OmniPath tend to score highly for at least one of their substrate's phosphosites (one-sided Wilcoxon rank sum test, $W = 1.0 \times 10^8$, $p < 1 \times 10^{-6}$) (Figure 2G, left). Furthermore, simply having a substrate site with a high functional score, indicating that the substrate is amenable to regulation by phosphorylation, can be predictive of a regulatory relationship (one-sided Wilcoxon rank sum test, $W = 6.4 \times 10^7$, $p < 1 \times 10^{-6}$) (Figure 2G, center). Linking these two metrics across all sites on the substrate via the DCG, we produced a score that could discriminate true regulatory relationships (one-sided Wilcoxon rank sum test, $W = 4.1 \times 10^7$, $p < 1 \times 10^{-6}$) (Figure 2G, right).

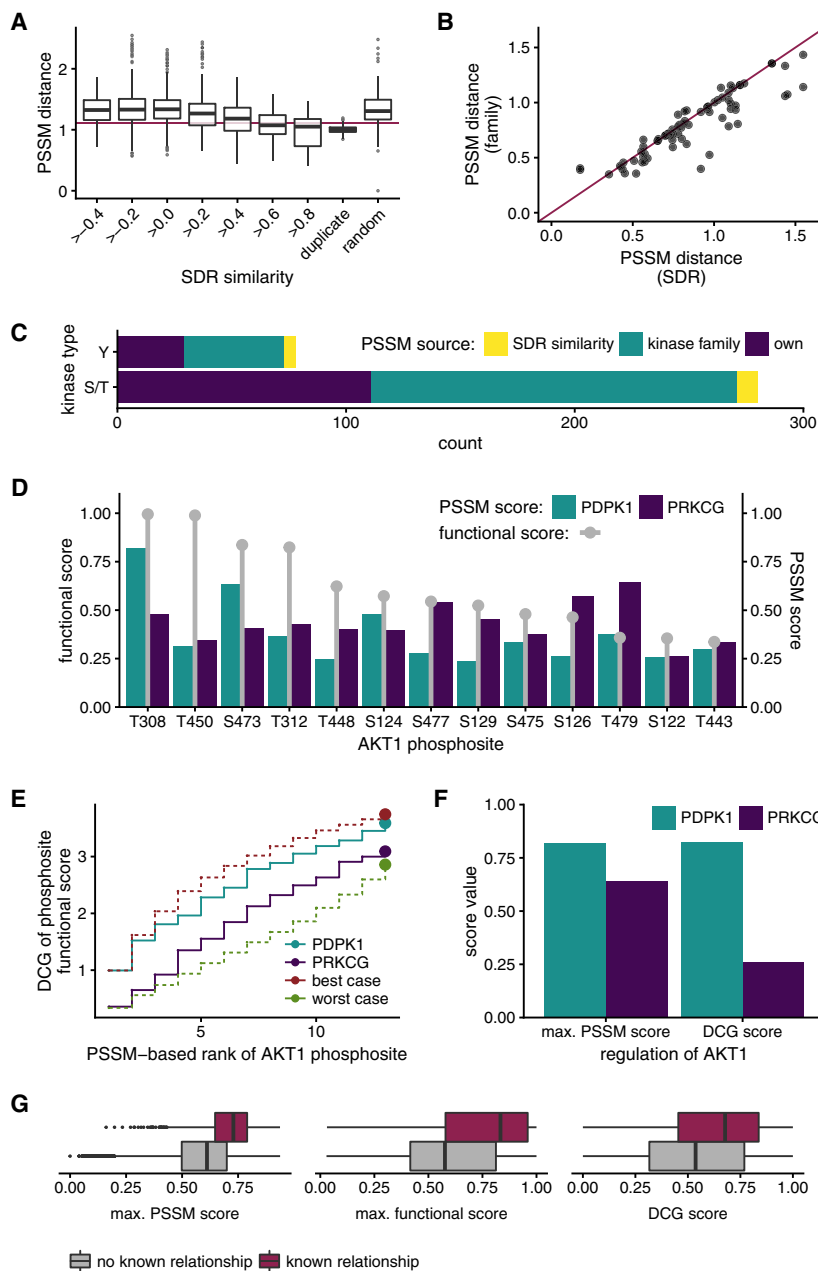


Figure 2. Kinase-Kinase Regulatory Relationships Can Be Predicted from Sequence Specificity and Phosphosite Functional Scores

(A) Similar kinase SDRs also indicate similar PSSMs. The red line indicates the 97.5th percentile of the distribution of distances between cross-validation PSSMs using different subsets of a kinase’s annotated substrates. At an SDR similarity of at least 0.8, over 50% of assigned PSSMs are less than this distance from their true values. (B) Assigning family-wise, composite PSSMs to unannotated kinases achieves similar, if not better, performance to SDR-based assignment. (C) Numbers of PSSMs by source (own annotations = 140, by family = 208, and by SDR similarity = 14). (D) PSSMs locate functional sites on substrates with differing performance. Here, the PSSM of PDPK1, a known regulator of AKT1, scores highly for sites with high functional scores, while that of PRKCG does not. (E) DCG quantifies the potential for a kinase to phosphorylate a putative substrate at its functional sites. (F) Although both PDPK1 and PRKCG have similar maximum PSSM scores for phosphorylating AKT1, only PDPK1 achieves a high DCG. (G) Maximum PSSM score, maximum substrate-site functional score, and DCG all discriminate regulatory relationships annotated in the OmniPath database from unannotated ones.

(Strumillo et al., 2019): sites within hotspots tend overwhelmingly to be activating (i.e., within the kinase activation loop) (one-sided Wilcoxon rank sum test, $W = 1.5 \times 10^3$, $p < 1 \times 10^{-6}$) (Figure 3A, first panel). When considering the sites’ positions within the domain, we found that most inhibitory sites are N-terminal (one-sided Wilcoxon rank sum test, $W = 1.5 \times 10^3$, $p = 0.045$) (Figure 3A, second panel). On the other hand, inhibitory sites tended to be more C-terminal in the overall protein, although the difference was not significant, (one-sided Wilcoxon rank sum test, $W = 7.5 \times 10^3$, $p = 0.38$) (Figure 3A, third panel). Lastly, we also observed that activating sites tend to be in more structured regions of the protein and inhibitory sites are more likely to be disordered,

although 50% of all inhibitory sites still were predicted to be in structured regions (one-sided Wilcoxon rank sum test, $W = 5.3 \times 10^3$, $p = 1.3 \times 10^{-4}$) (Figure 3A, fourth panel).

We then trained a predictor of phosphosite regulatory sign using these features (Table S1) via the Bayesian additive regression trees (BART) method. Cross-validation of the model showed consistently good performance, with a maximum mean Matthew’s correlation coefficient of 0.42 at a cutoff of 0.58 (posterior probabilities lower than the cutoff are declared to reflect inhibitory functionality), indicating overall good sign-classification performance (Figure 3B). Adjusting these posterior probabilities by the highest-performing cutoff provided us with a sign score for all phosphosites in our dataset, with negative scores indicating a prediction of

Protein Sequence and Structure Discriminate Phosphosites that Induce or Inhibit Kinase Activity

Phosphorylation events can lead to different regulatory outcomes for the substrate kinase, potentially inducing or inhibiting its enzymatic activity. Knowing these regulatory effects is essential to understanding the flow of information across complex networks of regulatory relationships. Thus, we sought to infer the “signs” (activating or inhibiting) of regulatory relationships from data.

To do so, we first evaluated how phosphorylation at a specific site is likely to affect a given kinase’s activity, according to annotations from the PhosphoSitePlus database (Hornbeck et al., 2015). We found particular discrimination for sites within phosphorylation hotspots of the protein-kinase domain

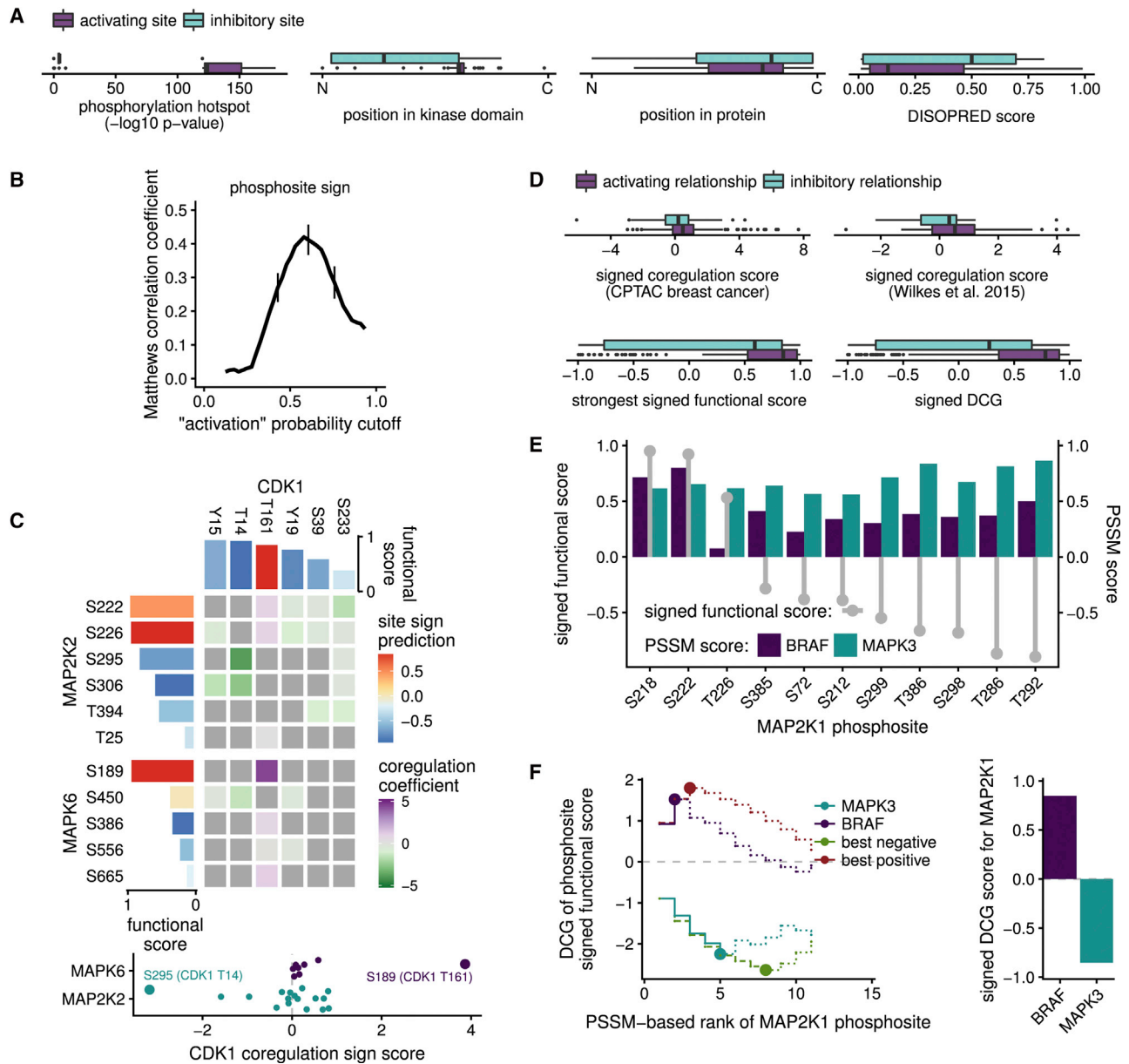


Figure 3. Evidence of Regulatory Sign (Activating versus Inhibiting) Can Be Uncovered in a Data-Driven Manner

(A) The regulatory sign of a single phosphosite, as annotated in PhosphoSitePlus, can be discriminated by using structural information: whether the site is in a phosphorylation hotspot, where the site is within the protein-kinase domain (N, N-terminal; C, C-terminal), the relative position of the site within the protein (N, N-terminal; C, C-terminal), and whether the site is in a disordered region.

(B) Matthews correlation coefficients for different posterior probability cutoffs for the predictor of phosphosite regulatory sign. The cutoff (above which a site or relationship would be declared to be "activating") that maximizes the coefficient discriminates best between inhibitory and activating sites or relationships. Error bars represent 95% confidence intervals.

(C) Modifying the phospho-coregulation score to account for predicted phosphosite sign and correlation sign can produce protein-level predictions of regulatory sign. Here, CDK1 is shown to have an activating relationship with MAPK6 via S189 and an inhibitory relationship with MAP2K2 via S295. Gray cells indicate missing or removed values.

(D) The signed variants of the coregulation score, functional score, and DCG, all discriminate between inhibitory and activating kinase-kinase regulatory relationships annotated in OmniPath.

(E) Accounting for predicted phosphosite sign can assess the propensity of a kinase to phosphorylate activating or inhibiting sites: BRAF's PSSM scores highly for activating sites on MAP2K1, while MAPK3 scores highly for inhibitory sites.

(F) A modified DCG for signed functional scores correctly assigns BRAF as an activator of MAP2K1 and MAPK3 as an inhibitor. Because there are more inhibitory sites on MAP2K1, a full DCG would be negative in most cases (dotted lines). Instead, we take the most extreme value visited by the sum (solid lines).

inhibition of activity and positive scores predicting activation (Table S2).

Kinase Regulatory Sign Can Be Inferred from Phosphosite Sign and Interaction Evidence

With phosphosite sign predictions in hand, we aimed to predict the signs of kinase-kinase regulatory interactions. Returning to the coregulation of functional phosphosites, we tested the consistency of the observed phosphoproteomic correlation with the sign of the phosphorylating kinase's regulatory site. If phosphorylation of an inhibitory site on a kinase is anticorrelated with that of an activating site on a putative substrate, then the evidence would suggest that the kinase positively regulates the substrate's activity. On the other hand, no direct regulation scenario would explain a positive correlation between these sites.

For example, cyclin-dependent kinase 1 (CDK1) shows strong evidence of negative coregulation with dual specificity mitogen-activated protein kinase kinase 2 (MAP2K2) at its site S295, reflecting its role in inhibiting MAP kinase kinases (Rossomando et al., 1994) (Figure 3C). CDK1 also activates MAPK6 (Tanguay et al., 2010) and, indeed, we find a strong positive correlation between two activating sites (CDK1 T161 and MAPK6 S189) on these kinases (Figure 3C). Overall, we found that the signed coregulation score was able to discriminate between activating and inhibitory kinase regulatory relationships, as annotated in OmniPath, in both phosphoproteomic datasets with activating relationship tending to have more positive coregulation score even though, in the case of the Wilkes study, the difference was not significant (one-sided Wilcoxon rank sum test, $W = 6.8 \times 10^2$, $p = 0.079$ [Wilkes et al., 2015] and $W = 1.1 \times 10^4$, $p = 0.0067$ [Mertins et al., 2016]) (Figure 3D, first and second panels).

We also adapted our DCG methodology after applying our sign predictions to the site functional scores. Thus, we now asked whether a kinase's PSSM tends to find relevant activating sites or inhibitory sites. For example, dual specificity MAP2K1 is activated by serine/threonine-protein kinase B-raf (BRAF) (Alessi et al., 1994; Macdonald et al., 1993; Papin et al., 1995) and is inhibited in negative feedback by its downstream substrate, MAPK3 (Eblen et al., 2004). We found that, indeed, B-raf has specificity toward MAP2K1's activating sites while MAPK3 is specific toward the inhibitory sites (Figure 3E). We then calculated a DCG on the signed functional scores, taking the most extreme value visited by the sum (Figure 3F). This method provides a positive value for BRAF and a negative value for MAPK3, as expected. Overall, both the signed functional score and the signed DCG score could discriminate well between activating and inhibitory relationships (one-sided Wilcoxon rank sum test, $W = 2.1 \times 10^4$, $p < 1 \times 10^{-6}$ [signed DCG] and $W = 2.3 \times 10^4$, $p < 1 \times 10^{-6}$ [signed functional score]) (Figure 3D, third and fourth panels). However, predictions for inhibitory relationships overall were less reliable.

A Supervised Learning Model Predicts a Global Network of Kinase Regulatory Relationships from Diverse Features

We combined the above evidence into two predictors via machine learning. The edge predictor predicts whether a kinase-ki-

nase regulatory relationship exists. The sign predictor predicts whether a given relationship induces or inhibits the substrate's kinase activity.

For training and validating the edge predictor, we retrieved from the OmniPath meta-database (Türei et al., 2016) a list of annotated relationships with at least two source databases supporting them, comprising 825 interactions in all. Because it is more difficult to prove the absence of a regulatory relationship, there is a lack of annotations for genuinely false relationships. We assumed that, in the space of all possible kinase-kinase interactions, regulatory relationships are rare. Therefore, a randomly selected pair of kinases is unlikely to show any regulatory relationship. We thus assessed the features described above for their predictive power on a validation set consisting of the annotated positive cases and random "negative" subsets of the remaining space of putative interactions.

Overall, each of the edge predictor features (Table S3) exhibited limited but measurable predictive power. We visualized this by the receiver operating characteristic (ROC) curve, comparing true-positive and false-positive rates as the score cutoff for declaring a regulatory relationship is lowered, and by similarly assessing precision and recall across cutoffs (Figure 4A). Maximum PSSM score performed the best, with a mean area under the ROC curve (AUC) of 0.742 ($\sigma = 0.007$, $n = 100$) (Figure S1A). The remaining features had mean AUC values of less than 0.7. We also noted that the precision decayed rapidly with lower cutoffs.

We then combined these features into the edge predictor with the BART method (Chipman et al., 2010) (Table S4). We first performed 3-fold cross-validation on the model 20 times using different random iterations of the training set (Figure 4A). The resulting models had a mean AUC of 0.884 ($\sigma = 0.009$, $n = 60$), representing a significant improvement over the individual features (Figure S1A).

We applied the same BART method to the regulatory sign features (Table S5) to produce the sign predictor (Table S4). We trained the model using regulatory signs annotated in OmniPath and evaluated it through cross-validation. Overall, performance was similar to the underlying site-level predictor described above, with a mean maximum Matthews correlation coefficient of 0.42; however, confidence intervals over the cross-validation were narrower for kinase-level predictions than they were for site-level predictions (Figure 4B). The maximum correlation occurred at a cutoff of 0.484 (i.e., the probability above which we would declare regulation to activate the substrate).

We next considered whether known, annotated relationships tend to rank highly among our edge predictions for each kinase. After building a new model for each kinase without using any of its annotated relationships in the training set, we found that 50% of kinases had a known regulatory relationship within the top 10 of our predictions (Figure 4C). The top ranks were significantly better than expected, based on random, per-kinase permutations of the scores (one-sided Wilcoxon rank sum test, regulator: $W = 5,818.5$, $p < 1 \times 10^{-6}$; substrate: $W = 7,385.5$, $p < 1 \times 10^{-6}$).

To further evaluate our model, we looked at how well it predicted interactions that were not included in the positive set due to being supported by only one source in OmniPath

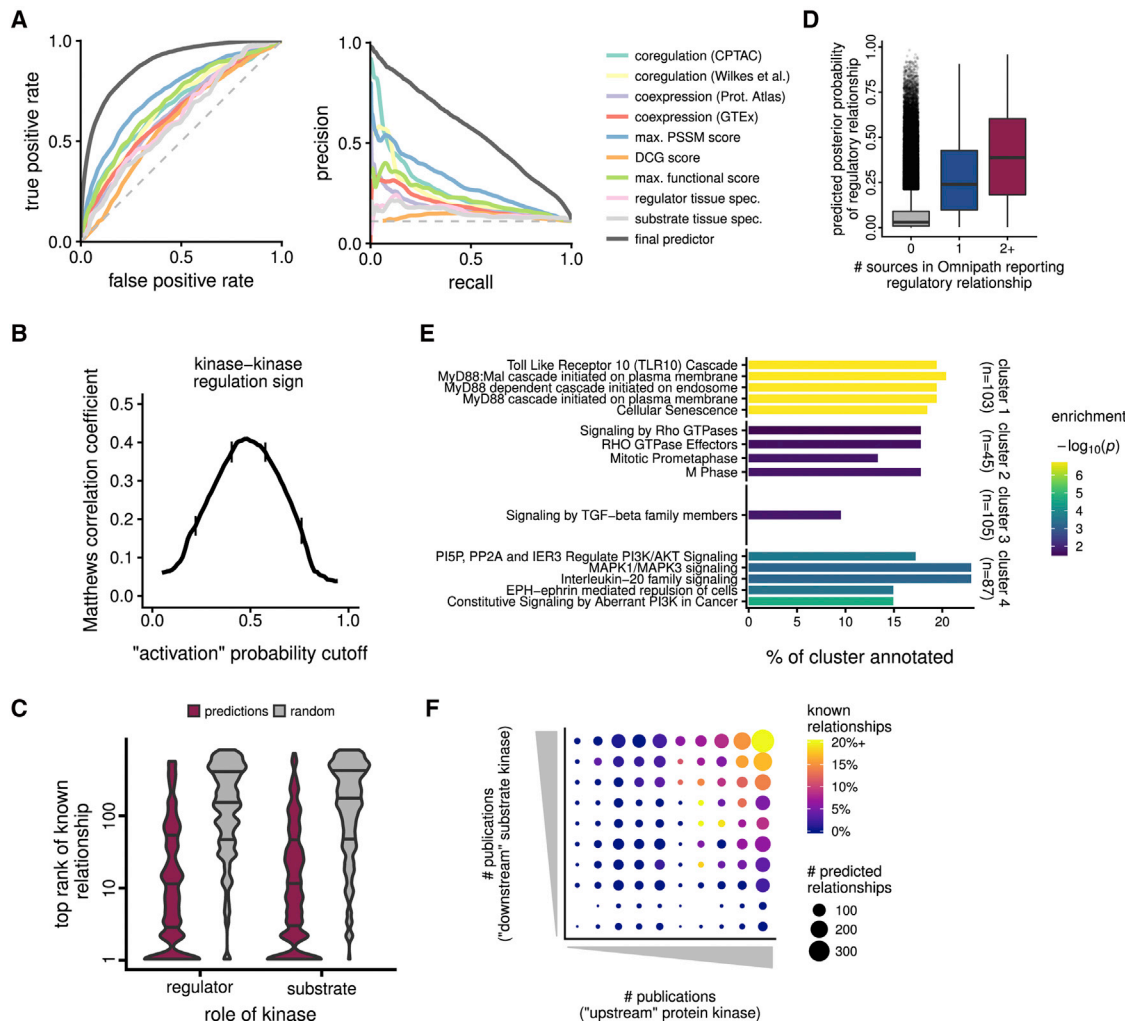


Figure 4. Combining Data-Driven Predictors of Kinase-Kinase Regulatory Relationships

(A) ROC and precision-recall curves of each feature and the final edge predictor. See also Figure S1A.

(B) Matthews correlation coefficients for different posterior probability cutoffs for the sign predictor. Error bars represent 95% confidence intervals.

(C) Annotated regulatory relationships for each kinase tend to rank highly among the predictions, when considering the kinase as either a regulator or a substrate. Lines indicate quartiles. 50% of kinases had a known regulatory relationship in the top ten predictions, which is significantly better than random expectation. See also Figure S2B.

(D) Previously annotated relationships supported by only one source in OmniPath score similarly to those supported by two or more sources (used in our training set), further validating our predictions.

(E) Clusters identified on the regulatory sub-network at a posterior probability cutoff of 0.5 are significantly enriched in annotations for unique sets of pathways. See also Figures S1B and S1C.

(F) The predicted network expands upon the annotated network, especially for understudied protein kinases. See also Figure S2A.

($n = 293$). These interactions had significantly higher prediction scores than unannotated regulatory interactions; however, they were generally lower than the high-confidence set (one-sided Wilcoxon rank sum test versus unannotated: $W = 6 \times 10^7$, $p < 1 \times 10^{-6}$, versus high-confidence set: $W = 87,073$, $p < 1 \times 10^{-6}$; Figure 4D).

We also noted several high-probability predictions that, while not being annotated in OmniPath, have direct or plausible support in the literature. For example, we predict receptor tyrosine protein kinase erbB-2 (ERBB2/HER2) to activate ephrin type-A receptor 2 (EPHA2) (edge probability 0.94 and sign probability 0.79). These two oncogenic kinases form a complex, and in a

mouse model of breast cancer, they appear to cooperate to promote tumor progression (Brantley-Sieders et al., 2008); however, no direct regulatory relationship has yet been described. We also predict that the closely related tyrosine protein kinases Fer (FER) and Fes/Fps (FES) activate HGF receptor (MET) with equal probabilities (edge probabilities 0.92 and sign probabilities 0.80). In fact, activation of MET by FER has previously been reported (Fan et al., 2016); however, this relationship is not annotated in OmniPath and thus was not present in our training and validation set. As a final example, we predict tyrosine protein kinase ABL1 to activate focal adhesion kinase 1 (PTK2/FAK1) (edge probability 0.92, sign probability 0.81). While, to our knowledge, no such

regulatory relationship has been described previously, FAK1 plays an important role in acute lymphoblastic leukemia characterized by constitutively active ABL1, and its phosphorylation was speculated to be “likely augmented by the direct action of activated ABL1 itself” (Churchman et al., 2016).

We next assessed the topology of regulatory relationships using a sub-network of high-confidence predictions (probability greater than 0.5), consisting of 340 kinases and 4,339 regulatory relationships (representing less than 2% of all possible relationships). We first applied a cluster-detection algorithm to an undirected variant of this network (retaining the higher probability relationship when two kinases were predicted to regulate each other, producing 3,716 undirected edges). Four clusters were detected (103, 45, 105, and 87 kinases, respectively; Table S6). This division of the network had a modularity of 0.325, which was significantly higher than expected given the modularity of randomized networks with the same degree distribution ($\mu = 0.197$, $\sigma = 0.00339$; $p < 0.001$; Figure S1B). To determine if these clusters reflected known biological associations, we tested each one for enrichment in pathway annotations from Reactome (Fabregat et al., 2017). Each cluster was enriched in annotations for at least one distinct Reactome pathway, indicating that the network successfully identified clusters of physiologically related kinases (Figure 4E; Table S7). We also assessed how related the pathways associated with each cluster were, using the average number of reactions between the proteins of two pathways as a proxy for relatedness. We found that the pathways associated with the same cluster were more closely related to each other than to those associated with other clusters ($p < 1 \times 10^{-6}$, $W = 5.8 \times 10^5$, Wilcoxon rank sum test; Figure S1C).

Because we set out to overcome the shortcomings inherent to literature-derived signaling pathway annotations, we checked for relationships between kinase connectivity on the high-confidence network and kinase publication counts (Figure 4F; kinase publication counts were retrieved from, Invergo and Beltrao, 2018). Interactions between kinases in the top three publication-count deciles (more than 95 publications) accounted for only 31% of the network. Conversely, 589 regulatory relationships were predicted between pairs of kinases in the bottom 50% of publication counts (fewer than 40 publications each).

Overall, only 7% of the relationships in the high-confidence network are annotated in databases. This portion increases as a higher probability threshold is applied to the network (Table S8). Although the number of previously annotated interactions is dwarfed by novel predictions, a significant proportion of this can be accounted for by the relative sparsity of annotated relationships for understudied kinases. Restricting the network to highly studied kinases largely resolves this (Figure S2A). However, this can also be explained in part by a persistence of the influence of better annotation for well-studied kinases in our predictions, as can be seen in a significant correlation between publication count and top prediction-rank of known relationships (Figure S2B; Spearman’s rank correlation, as regulator: $\rho = -0.34$, $p < 1 \times 10^{-6}$; as substrate: $\rho = -0.29$, $p < 1 \times 10^{-6}$).

The Trained Model Can Reconstruct Known Signaling Pathways

We next investigated whether our data-driven, signed kinase-kinase regulatory predictions were able to reconstruct known

pathways. For each kinase that we include, we generated a new model for which all regulatory relationships including the kinase were left out of the training set. These kinase-specific models were then used to predict the kinase’s regulatory substrates and the signs of the interactions. To these we applied an edge probability cutoff of 0.5 and a sign cutoff of 0.5. We started by choosing well-studied kinases that are functionally related to AKT1 (Figure 5A). Between these kinases, we successfully recovered all but one annotated relationship, the regulation of ribosomal protein S6 kinase beta-1 (RPS6KB1) by AKT1 and PDPK1. Six predicted relationships are not present in database annotations. Sign predictions generally fail on a per-substrate basis. For example, we predict all regulations of serine/threonine-protein kinase mTOR (MTOR) to be inhibitory, while those that have been annotated are activating. Our predictions also perform well when considering MAPK signaling, recovering all but three previously annotated edge, and with two erroneous predictions of an annotated sign (Figure 5B).

If we begin to include other paralogs of these kinases, which tend to be less well studied, we quickly accumulate predictions for previously undescribed relationships. For example, we predict many modes of inter-regulation between S6 kinases and glycogen synthase kinases. On the other hand, we fail to predict several regulatory relationships involving RAC protein kinases AKT2 and AKT3 (Figure 5C). Expanding the MAPK signaling network is more successful, with the core signaling events being recovered between RAFs, MAP2Ks, and MAPKs, including correct sign prediction, while also filling in the network of interactions for the less well-studied A-Raf (ARAF) and MAP2K3 (Figure 5D). Both of these examples demonstrate that the predicted networks quickly become difficult to assess visually when more than a few kinases are included, particularly those with fewer annotations, because of the numbers of unvalidated predictions. However, extrapolating from the overall performance and the success on smaller networks, our results suggest that this complexity is inherent to kinase signaling networks.

Independent Experimental Data Support Predicted Regulatory Relationships

We next investigated whether our predictions are reflected in kinase-target relationships identified in large-scale phosphoproteomics experiments that were not used for training. First, we employed two recently published datasets. In one, Sugiyama et al. (Sugiyama et al., 2019) have identified *in vitro* substrate phosphosites for 354 human kinases. In the second study, *in vivo* phosphosites that are directly or indirectly “downstream” of 103 kinases were determined by phosphoproteomic experiments after chemical inhibition of kinase activity (Hijazi et al., 2020). Together, these two datasets define kinase-substrate phosphosites relationships, from which we selected kinase-kinase phosphorylation relationships that we reasoned should be enriched in regulatory interactions. We note however that these studies identify whether a kinase could be responsible (directly or indirectly) for the phosphorylation of another but not necessarily whether such phosphorylation is regulatory.

We looked at probability scores assigned to relationships that were corroborated by these experiments. Relationships included in our validation set were discarded. In both cases, we observed that the probability score derived from our model was significantly higher for these experimentally identified

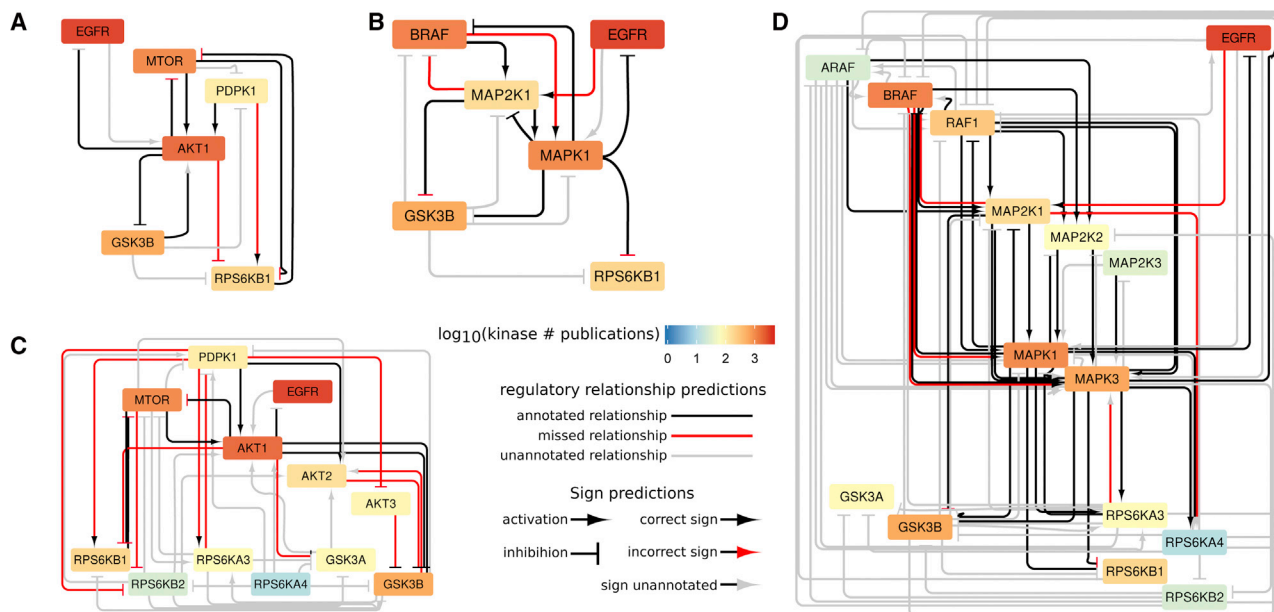


Figure 5. Our Data-Driven Predictor Reconstructs Known Signaling Pathways “from Scratch”

Each kinase’s outgoing edges and signs were predicted from models trained after leaving the kinase out from the training set. A probability cutoff of 0.5 and “activating” sign cutoff of 0.5 were used. Black edges are relationships annotated in OmniPath that were correctly recovered. Red edges are annotated relationships that were not predicted. Gray edges are unvalidated predictions. Arrowheads indicate the predicted regulatory sign: arrows indicate activation and bars indicate inhibition. Black arrows are correctly predicted, red arrows are incorrectly predicted, and gray arrows are unvalidated. Node colors indicate the number of publications associated with the kinase. Kinases with fewer associated publications (“cool” blue/green colors) tend to have more unvalidated edges. (A) A reconstruction of AKT1 signaling using only well-studied kinases largely recovers the known information flow of the pathway. (B) Similar performance is seen in reconstructing MAP kinase signaling. (C) Including lesser-studied kinases in the AKT signaling analysis greatly increases the number of unannotated or missed relationships while also predicting complex modes of regulatory feedback. (D) Expanding the MAP kinase signaling pathway to include more paralogs captures a highly interconnected core of previously annotated relationships while adding numerous unvalidated relationships between lesser-studied kinases.

kinase-kinase relationships than expected by chance (*in vitro*: median probability of 0.075 ($W = 2.63 \times 10^8$, $p < 1 \times 10^{-6}$); *in vivo*: median probability of 0.12 ($W = 1.57 \times 10^8$, $p < 1 \times 10^{-6}$); background: median probability of 0.030) (Figure 6A). In addition, a subset of these experimental target sites was found at positions with a higher regulatory potential based on the phosphosite functional score (Ochoa et al., 2020). When filtering the *in vitro* and *in vivo* kinase-kinase interactions by the phosphosite functional score cutoff of 0.5 we observe an increase in the median probability from our model (Figure 6A), in particular for the *in vivo* kinase-kinase target set (Figure 6A, “*in vivo* fnc”). Both the *in vitro* set ($W = 1.03 \times 10^8$, $p < 1 \times 10^{-6}$) and *in vivo* set ($W = 8.81 \times 10^7$, $p < 1 \times 10^{-6}$) filtered by the phosphosite functional score have a significantly higher edge probability than the background set. Furthermore, predicted network edges corroborated by either study had a higher probability of being included in our validation set with Fisher’s exact test (OR = 3.98 and $p = 1.6 \times 10^{-4}$ for the *in vivo* study and OR = 4.51 and $p < 1 \times 10^{-6}$ for the *in vitro* study). We provide in Table S9 the list of kinase-kinase regulatory relationships that have a high predicted score from our model having also *in vivo* or *in vitro* supporting evidence. This includes 3 cases of unannotated kinase-kinase relationships with support from both the *in vitro* and *in vivo* experiments.

Finally, as an application of our predictions, we tested whether phosphorylation perturbation data could be used to discover novel pathways within the inferred network. To this end, we measured changes in phosphorylation after treatment of human cells (Kasumi-1 cell line) with MEK and PI3K inhibitors via phosphoproteomics (Table S10). A total of 9,183 phosphopeptides were quantified for MEK and PI3Ki and control condition, with 66 and 112 phosphosites identified as significantly downregulated after inhibition. Of measured phosphosites, 650 had a known upstream kinase included in our high-probability network leading to 1,019 kinase-substrate interactions being added to the network. After omitting sites known to be direct substrates of MEK and PI3K as well as sites without any known kinase, 6 and 11 downregulated phosphosites were considered for subsequent analysis. We then asked if the inhibition of MEK (MAP2K1 and MAP2K2) and PI3K kinases could be linked to the downregulated phosphosites via connections predicted by our model. A probability cutoff of 0.5 was used to retain a network of highly probable edges and owing to the fact that PI3K is a lipid kinase, PI3K was linked to the network via its known substrate kinases. The regulated phosphosites were added to the predicted kinase network based on prior knowledge, and the distances between the inhibited kinases and the downregulated phosphosites were calculated as the sum of weights across the shortest weighted path on the predicted network. We found that

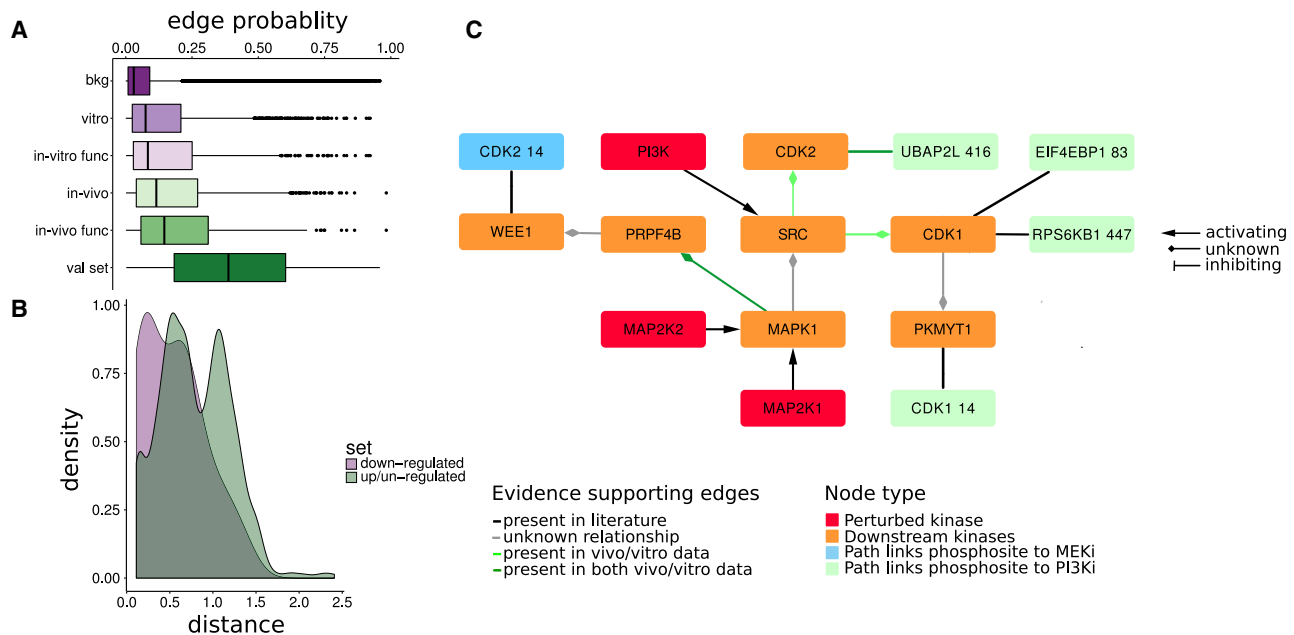


Figure 6. Network Validation with Independent Experimental Data

(A) Edges corroborated by either *in vitro* or *in vivo* datasets had higher probability, 0.075 ($n = 1,596$) and 0.12 ($n = 856$), respectively, compared with otherwise unsupported edges, 0.030 ($n = 251,159$). Edges filtered by functional score had yet higher probability assigned to them; 0.083 ($n = 602$) and 0.15 ($n = 456$). (B) Downregulated phosphosites tended to be closer in terms of weighted shortest path to perturbed kinases. (C) Connecting downregulated phosphosites to perturbed kinases by traversing through the shortest weighted path in the network yielded two predicted interactions corroborated by the *in vitro* data (SRC and CDK1/2) and one predicted interaction supported by both the *in vitro* and *in vivo* datasets (MAPK1 and PRPF4B). Known edges were derived from the OmniPath validation set apart from the edge between PI3K and SRC, which was derived from KEGG.

downregulated phosphosites were closer to the inhibited kinases than all other quantified phosphosites ($W = 1.4 \times 10^4$, $p = 0.0028$) (Figure 6B). We then selected paths connecting the inhibited kinases with downregulated phosphosites via edges that had support from at least one of the kinase-target experiments described above (Figure 6C). In this way, we identified putative kinase regulatory interactions that have a high-probability score (>0.5), are supported by *in vitro* or *in vivo* kinase-target experiments, and help to explain downregulation of phosphosites from our kinase inhibition experiments. This includes 2 related relationships (Figure 6C) between SRC and two CDKs (CDK1 and CDK2) and the predicted regulation of PRPF4B by MAPK1. The latter was supported by both the *in vitro* and *in vivo* studies.

DISCUSSION

The task of experimentally testing all possible kinase-kinase relationships in order to produce a complete regulatory network is daunting. We have thus taken a data-driven supervised machine learning-based approach to predict these regulatory relationships. Although we do not suggest that these predictions can replace established methods for confirming regulatory relationships, they can nevertheless be used to reduce the vast space of possible relationships under consideration in order to form credible hypotheses and to prioritize experiments, particularly for understudied kinases.

Previous efforts to produce kinome-scale inferences of regulatory relationships have depended on scaffolding data-driven

predictions to existing protein networks. For example, Rudolph et al. (2016) derived signaling pathways through a network diffusion technique with phosphoproteomic data on a literature-derived protein-protein interaction network. However, such analyses are strongly impacted by the incompleteness of the existing networks and the overrepresentation of well-studied kinases therein (Gillis et al., 2014; Invergo and Beltrao, 2018; Rolland et al., 2014). To our knowledge, there has only been one other attempt to predict kinase regulatory sign (Hernandez et al., 2010). The authors inferred signs from quantitative phosphoproteomic data on a literature-derived kinase network, in which the method in part depended upon the connectivity of the kinases on this network. However, missing or erroneous annotated relationships could have major impacts on the results. Our supervised machine learning approach can make predictions for kinases with no previously known information, in this way improving the coverage in our predicted network. We only retain a literature influence, which might affect the generality of our model, from using annotated substrates in the construction of our kinase specificity models and from the construction of the training set. The former can be resolved with high-throughput methods to measure kinase specificity profiles (see, e.g., Imamura et al., 2014, Sugiyama et al., 2019). The latter, which could omit highly specialized modes of regulation, can be improved as more relationships are experimentally validated.

Many factors can affect the nature of a kinase-kinase regulatory relationship and each such relationship will be unique, owing to the particular properties of the kinases involved. Thus, making generalized predictions that apply to all of them

is inherently difficult. Nevertheless, some features are fundamental, such as regulation by phosphorylation, and identification of how much each of such features contribute to the specificity in kinase-substrate interaction will be key in predicting the regulatory roles of phosphorylation sites. To this end, the performance of the predictions via identifying patterns of phosphorylation will improve with more data. Given the importance of PSSMs in our results, there is a clear need for producing robust PSSMs for every kinase in order to prune indirect regulatory effects. As for correlative methods on phosphoproteomic data, many conditions are needed to confidently discriminate the phosphoregulation of over 500 kinases. Importantly, large-scale phosphoproteomics experiments are needed across a more diverse array of tissues or cell lines to properly capture the activities of more tissue-specific kinases. Because we only used data from experiments using the breast cancer cell line MCF7, many kinases were not represented in the phosphoproteomic data. Furthermore, the use of data derived from cancer cell lines might introduce errors in the resulting network since cancer initiation and progression disrupt intracellular signaling (Sever and Brugge, 2015).

We assumed in the construction of our predictor that the true network is sparse, and indeed we assign to 75% of all possible relationships posterior probabilities of less than 0.09, far below any probability cutoff that we considered. Nevertheless, even at stringent cutoffs, isolating a sub-network of more than a few kinases produces a denser topology of regulatory relationships than is typically considered for kinase signaling. It is possible that this is an artifact of not considering cellular context (e.g., protein expression or cellular localization). There is also an unavoidable accumulation of false-positives as more predictions are considered. Despite these caveats, our results suggest that the kinase regulatory network is richer in feedback and cross-module regulation than expected based on the current view of kinase pathways. Further developments in experimental approaches for hypothesis-free kinase regulatory network reconstruction are needed to confirm the predicted modularity and density of regulatory relationships in kinase signaling networks.

Key Changes Prompted by Reviewer Comments

In response to reviewer comments, we have added validation analyses on the external *in vitro* (Sugiyama et al., 2019) and *in vivo* (Hijazi et al., 2020) datasets, as well as validation on a newly generated phosphoproteomic dataset. This resulted in the addition of Figure 6. We also added statistical tests to support the assessment of each feature's discriminatory power. We updated the methodology used to generate Figure 5 in order to assure that each kinase was removed from the training set before predicting its regulatory relationships. Finally, we added Table S8 to better illustrate the network at different score cutoffs. For context, the complete transparent peer review record is included within the Supplemental Information.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE
- LEAD CONTACT AND MATERIALS AVAILABILITY
- EXPERIMENTAL MODEL AND SUBJECT DETAILS
- METHOD DETAILS
 - Data
 - Protein Kinase Specificity Models
 - Scoring Phosphosites with PSSMs
 - Phosphosite Functional Scores
 - Linking PSSMs to Phosphosite Functional Scores
 - Coexpression and Tissue Specificity
 - Phospho-Coregulation
 - Prediction of Regulatory Relationships
 - Network Clustering and Pathway Enrichment
 - Pathway-Annotation Distances
 - Network Modularity
 - Prediction of Phosphosite Functional Sign
 - Prediction of Regulatory Sign
 - Signed Discounted Cumulative Gain
 - Signed Coregulation Score
 - Training and Validation of the Sign Predictor
 - Kinase Inhibitor Experiments
 - Identification of Down-regulated Phosphosites
 - Kinase-Substrate Shortest Paths
- QUANTIFICATION AND STATISTICAL ANALYSES
- DATA AND CODE AVAILABILITY

SUPPLEMENTAL INFORMATION

Supplemental Information can be found online at <https://doi.org/10.1016/j.cels.2020.04.005>.

ACKNOWLEDGMENTS

B.M.I. received funding from F. Hoffmann-La Roche Ltd.

AUTHOR CONTRIBUTIONS

P.B., E.P., and B.M.I. conceived of the study; B.M.I. and B.P. performed the analyses and wrote the manuscript; N.A. and P.C. conceived the phosphoproteomic experiments that were carried out by N.A. and M.H. D.B. and G.G. performed additional analyses.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: July 15, 2019
Revised: January 24, 2020
Accepted: April 20, 2020
Published: May 20, 2020

REFERENCES

- Acosta-Jaquez, H.A., Keller, J.A., Foster, K.G., Ekim, B., Soliman, G.A., Feener, E.P., Ballif, B.A., andingar, D.C. (2009). Site-specific mTOR phosphorylation promotes mTORC1-mediated signaling and cell growth. *Mol. Cell. Biol.* 29, 4308–4324.
- Alessi, D.R., Saito, Y., Campbell, D.G., Cohen, P., Sihanandam, G., Rapp, U., Ashworth, A., Marshall, C.J., and Cowley, S. (1994). Identification of the sites in MAP kinase kinase-1 phosphorylated by p74raf-1. *EMBO J.* 13, 1610–1619.
- Babur, Ö., Ngo, A.T.P., Rigg, R.A., Pang, J., Rub, Z.T., Buchanan, A.E., Mitrugno, A., David, L.L., McCarty, O.J.T., Demir, E., Aslan, J.E., et al. (2018). Platelet procoagulant phenotype is modulated by a p38-MK2 axis

- that regulates RTN4/Nogo proximal to the endoplasmic reticulum: utility of pathway analysis. *Am. J. Physiol. Cell Physiol.* **374**, C603–C615.
- Basson, M.A. (2012). Signaling in cell differentiation and morphogenesis. *Cold Spring Harb. Perspect. Biol.* **4**, a008151.
- Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. B* **57**, 289–300.
- Blom, N., Sicheritz-Pontén, T., Gupta, R., Gammeltoft, S., and Brunak, S. (2004). Prediction of post-translational glycosylation and phosphorylation of proteins from the amino acid sequence. *Proteomics* **4**, 1633–1649.
- Blondel, V.D., Guillaume, J.-L., Lambiotte, R., and Lefebvre, E. (2008). Fast unfolding of communities in large networks. *J. Stat. Mech.* **2008**, 10008.
- Bolstad, B.M., Irizarry, R.A., Astrand, M., and Speed, T.P. (2003). A comparison of normalization methods for high density oligonucleotide array data based on variance and bias. *Bioinformatics* **19**, 185–193.
- Bradley, D., and Beltrao, P. (2019). Evolution of protein kinase substrate recognition at the active site. *PLoS Biol* **17**, e3000341.
- Bradley, D., Vieitez, C., Rajeeve, V., Cutillas, P.R., and Beltrao, P. (2018). Global analysis of specificity determinants in eukaryotic protein kinases. [bioRxiv https://www.biorxiv.org/content/10.1101/195115v2](https://www.biorxiv.org/content/10.1101/195115v2).
- Brantley-Sieders, D.M., Zhuang, G., Hicks, D., Fang, W.B., Hwang, Y., Cates, J.M.M., Coffman, K., Jackson, D., Bruckheimer, E., Muraoka-Cook, R.S., et al. (2008). The receptor tyrosine kinase EphA2 promotes mammary adenocarcinoma tumorigenesis and metastatic progression in mice by amplifying ErbB2 signaling. *J. Clin. Invest.* **118**, 64–78.
- Cheng, G., Ye, Z.S., and Baltimore, D. (1994). Binding of Bruton's tyrosine kinase to Fyn, Lyn, or Hck through a Src homology 3 domain-mediated interaction. *Proc. Natl. Acad. Sci. USA* **91**, 8152–8155.
- Chipman, H.A., George, E.I., and McCulloch, R.E. (2010). BART: bayesian additive regression trees. *Ann. Appl. Stat.* **4**, 266–298.
- Churchman, M.L., Evans, K., Richmond, J., Robbins, A., Jones, L., Shapiro, I.M., Pachter, J.A., Weaver, D.T., Houghton, P.J., Smith, M.A., et al. (2016). Synergism of FAK and tyrosine kinase inhibition in Ph⁺ B-ALL. *JCI Insight* **1**, e86082.
- Clauset, A., Newman, M.E.J., and Moore, C. (2004). Finding community structure in very large networks. *Stat Nonlin Soft Matter Phys* **70**, 066111.
- Csárdi, G., and Nepusz, T. (2006). The igraph software package for complex network research. *InterJournal Complex Systems* **1695**.
- Dinkel, H., Chica, C., Via, A., Gould, C.M., Jensen, L.J., Gibson, T.J., and Diella, F. (2011). Phospho.ELM: a database of phosphorylation sites—update 2011. *Nucleic Acids Res.* **39**, D261–D267.
- Eblen, S.T., Slack-Davis, J.K., Tarcsafalvi, A., Parsons, J.T., Weber, M.J., and Catling, A.D. (2004). Mitogen-activated protein kinase feedback phosphorylation regulates MEK1 complex formation and activation during cellular adhesion. *Mol. Cell Biol.* **24**, 2308–2317.
- Ellis, J.J., and Kobe, B. (2011). Predicting protein kinase specificity: predikin update and performance in the DREAM4 challenge. *PLoS One* **6**, e21169.
- Fabregat, A., Jupe, S., Matthews, L., Sidiropoulos, K., Gillespie, M., Garapati, P., Haw, R., Jassal, B., Korninger, F., May, B., et al. (2017). The Reactome pathway knowledgebase. *Nucleic Acids Res.* **46**, D649–D655.
- Fan, G., Zhang, S., Gao, Y., Greer, P.A., and Tonks, N.K. (2016). HGF-independent regulation of MET and GAB1 by nonreceptor tyrosine kinase fer potentiates metastasis in ovarian cancer. *Genes Dev.* **30**, 1542–1557.
- Gillis, J., Ballouz, S., and Pavlidis, P. (2014). Bias tradeoffs in the creation and analysis of protein–protein interaction networks. *J. Proteomics.* **170**, 44–54.
- GTEx Consortium (2013). The genotype-tissue expression (GTEx) project. *Nat Genet.* **45**, 580–585.
- Henikoff, S., and Henikoff, J.G. (1994). Position-based sequence weights. *J. Mol. Biol.* **243**, 574–578.
- Henikoff, J.G., and Henikoff, S. (1996). Using substitution probabilities to improve position-specific scoring matrices. *Comput. Appl. Biosci.* **12**, 135–143.
- Hernandez, M., Lachmann, A., Zhao, S., Xiao, K., and Ma'ayan, A. (2010). Inferring the sign of kinase-substrate interactions by combining quantitative phosphoproteomics with a literature-based mammalian kinome network. *Proc IEEE Int Symp Bioinformatics Bioeng.* **2010**, 180–184.
- Hijazi, M., Smith, R., Rajeeve, V., Bessant, C., and Cutillas, P.R. (2020). Reconstructing kinase network topologies from phosphoproteomics data reveals cancer-associated rewiring. *Nat. Biotechnol.* **38**, 493–502.
- Hill, S.M., Heiser, L.M., Cokelaer, T., Unger, M., Nesser, N.K., Carlin, D.E., Zhang, Y., Sokolov, A., Paull, E.O., Wong, C.K., et al. (2016). Inferring causal molecular networks: empirical assessment through a community-based effort. *Nat. Methods* **13**, 310–318.
- Hill, S.M., Lu, Y., Molina, J., Heiser, L.M., Spellman, P.T., Speed, T.P., Gray, J.W., Mills, G.B., and Mukherjee, S. (2012). Bayesian inference of signaling network topology in a cancer cell line. *Bioinformatics* **28**, 2804–2810.
- Hill, S.M., Nesser, N.K., Johnson-Camacho, K., Jeffress, M., Johnson, A., Boniface, C., Spencer, S.E.F., Lu, Y., Heiser, L.M., Lawrence, Y., et al. (2017). Context specificity in causal signaling networks revealed by phosphoprotein profiling. *Cell Syst.* **4**, 73–83.e10.
- Hornbeck, P.V., Zhang, B., Murray, B., Kornhauser, J.M., Latham, V., and Skrzypek, E. (2015). PhosphoSitePlus, 2014: mutations, PTMs and recalibrations. *Nucleic Acids Res.* **43**, D512–D520.
- Imamura, H., Sugiyama, N., Wakabayashi, M., and Ishihama, Y. (2014). Large-scale identification of phosphorylation sites for profiling protein kinase selectivity. *J. Proteome Res.* **13**, 3410–3419.
- Invergo, B.M., and Beltrao, P. (2018). Reconstructing phosphorylation signaling networks from quantitative phosphoproteomic data. *Essays Biochem.* **62**, 525–534.
- Järvelin, K., and Kekäläinen, J. (2002). Cumulated gain-based evaluation of IR techniques. *ACM Trans. Inf. Syst.* **20**, 422–446.
- Kanehisa, M. (2019). Toward understanding the origin and evolution of cellular organisms. *Protein Sci.* **28**, 1947–1951.
- Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y., and Morishima, K. (2017). KEGG: new perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–D361.
- Kanehisa, M., and Goto, S. (2000). KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **28**, 27–30.
- Kanehisa, M., Sato, Y., Furumichi, M., Morishima, K., and Tanabe, M. (2019). New approach for understanding genome variations in KEGG. *Nucleic Acids Res.* **47**, D590–D595.
- Kapelner, A., and Bleich, J. (2015). Prediction with missing data via bayesian additive regression trees. *Can. J. Statistics* **43**, 224–239.
- Kapelner, A., and Bleich, J. (2016). bartMachine: machine learning with Bayesian additive regression trees. *J. Stat. Soft.* **70**, 1–40.
- Kholodenko, B.N., Kiyatkin, A., Bruggeman, F.J., Sontag, E., Westerhoff, H.V., and Hoek, J.B. (2002). Untangling the wires: a strategy to trace functional interactions in signaling and gene networks. *Proc. Natl. Acad. Sci. USA* **99**, 12841–12846.
- Köksal, A.S., Beck, K., Cronin, D.R., McKenna, A., Camp, N.D., Srivastava, S., MacGilvray, M.E., Bodik, R., Wolf-Yadlin, A., Fraenkel, E., et al. (2018). Synthesizing signaling pathways from temporal Phosphoproteomic data. *Cell Rep.* **24**, 3607–3618.
- Linding, R., Jensen, L.J., Ostheimer, G.J., van Vugt, M.A.T.M., Jørgensen, C., Miron, I.M., Diella, F., Colwill, K., Taylor, L., Elder, K., et al. (2007). Systematic discovery of in vivo phosphorylation networks. *Cell* **129**, 1415–1426.
- Macdonald, S.G., Crews, C.M., Wu, L., Driller, J., Clark, R., Erikson, R.L., and McCormick, F. (1993). Reconstitution of the Raf-1-MEK-ERK signal transduction pathway in vitro. *Mol. Cell. Biol.* **13**, 6615–6620.
- Manning, G., Whyte, D.B., Martinez, R., Hunter, T., and Sudarsanam, S. (2002). The protein kinase complement of the human genome. *Science* **298**, 1912–1934.
- Mérieu, K., Jacquot, S., Zeniou, M., Pannetier, S., Sassone-Corsi, P., and Hanauer, A. (2000). Activation of RSK by UV-light: phosphorylation dynamics and involvement of the MAPK pathway. *Oncogene* **19**, 4221–4229.

- Mertins, P., Mani, D.R., Ruggles, K.V., Gillette, M.A., Clauser, K.R., Wang, P., Wang, X., Qiao, J.W., Cao, S., Petralia, F., et al. (2016). Proteogenomics connects somatic mutations to signalling in breast cancer. *Nature* *534*, 55–62.
- Oates, C.J., Dondelinger, F., Bayani, N., Korkola, J., Gray, J.W., and Mukherjee, S. (2014). Causal network inference using biochemical kinetics. *Bioinformatics* *30*, i468–i474.
- Oates, C.J., and Mukherjee, S. (2012). Network inference and biological dynamics. *Ann. Appl. Stat.* *6*, 1209–1235.
- Obenauer, J.C., Cantley, L.C., and Yaffe, M.B. (2003). Scansite 2.0: proteome-wide prediction of cell signaling interactions using short sequence motifs. *Nucleic Acids Res.* *31*, 3635–3641.
- Ochoa, D., Jarnuczak, A.F., Viéitez, C., Gehre, M., Soucheray, M., Mateus, A., Kleefeldt, A.A., Hill, A., Garcia-Alonso, L., Stein, F., et al. (2020). The functional landscape of the human phosphoproteome. *Nat. Biotechnol.* *38*, 365–373.
- Papathodorou, I., Fonseca, N.A., Keays, M., Tang, Y.A., Barrera, E., Bazant, W., Burke, M., Füllgrabe, A., Fuentes, A.M.-P., George, N., et al. (2018). Expression Atlas: gene and protein expression across multiple studies and organisms. *Nucleic Acids Res.* *46*, D246–D251.
- Papin, C., Eychène, A., Brunet, A., Pagès, G., Pouyssegur, J., Calothy, G., and Barnier, J.V. (1995). B-Raf protein isoforms interact with and phosphorylate Mek-1 on serine residues 218 and 222. *Oncogene* *10*, 1647–1651.
- Park, H., Wahl, M.I., Afar, D.E., Turck, C.W., Rawlings, D.J., Tam, C., Scharenberg, A.M., Kinet, J.P., and Witte, O.N. (1996). Regulation of Btk function by a major autophosphorylation site within the SH3 domain. *Immunity* *4*, 515–525.
- Perfetto, L., Briganti, L., Calderone, A., Cerquone Perpetuini, A., Iannuccelli, M., Langone, F., Licata, L., Marinkovic, M., Mattioni, A., Pavlidou, T., et al. (2016). SIGNOR: a database of causal relationships between biological entities. *Nucleic Acids Res.* *44*, D548–D554.
- R Core Team (2016). R: a language and environment for statistical computing (R Foundation for Statistical Computing).
- Rawlings, D.J., Scharenberg, A.M., Park, H., Wahl, M.I., Lin, S., Kato, R.M., Fluckiger, A.C., Witte, O.N., and Kinet, J.P. (1996). Activation of BTK by a phosphorylation mechanism initiated by SRC family kinases. *Science* *271*, 822–825.
- Rhind, N., and Russell, P. (2012). Signaling pathways that regulate cell division. *Cold Spring Harb. Perspect. Biol.* *4*, a005942.
- Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., and Smyth, G.K. (2015). limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* *43*, e47.
- Rolland, T., Taşan, M., Charlotteaux, B., Pevzner, S.J., Zhong, Q., Sahni, N., Yi, S., Lemmens, I., Fontanillo, C., Mosca, R., et al. (2014). A proteome-scale map of the human interactome network. *Cell* *159*, 1212–1226.
- Rossomando, A.J., Dent, P., Sturgill, T.W., and Marshak, D.R. (1994). Mitogen-activated protein kinase kinase 1 (MKK1) is negatively regulated by threonine phosphorylation. *Mol. Cell. Biol.* *14*, 1594–1602.
- Rudolph, J.D., de Graauw, M., van de Water, B., Geiger, T., and Sharan, R. (2016). Elucidation of signaling pathways from large-scale Phosphoproteomic data using protein interaction networks. *Cell Syst.* *3*, 585–593.e3.
- Sever, R., and Brugge, J.S. (2015). Signal Transduction in Cancer. *Cold Spring Harb. Perspect. Med.* *5*.
- Sing, T., Sander, O., Beerenwinkel, N., and Lengauer, T. (2005). ROCr: visualizing classifier performance in R. *Bioinformatics* *21*, 3940–3941.
- Smith, J.A., Poteet-Smith, C.E., Malarkey, K., and Sturgill, T.W. (1999). Identification of an extracellular signal-regulated kinase (ERK) docking site in ribosomal S6 kinase, a sequence critical for activation by ERK in vivo. *J. Biol. Chem.* *274*, 2893–2898.
- Strumillo, M.J., Oplova, M., Vieitez, C., Ochoa, D., Shahraz, M., Busby, B.P., Sopko, R., Studer, R.A., Perrimon, N., Panse, V.G., and Beltrao, P. (2019). Conserved phosphorylation hotspots in eukaryotic protein domain families. *Nat. Commun.* *10*, 1977.
- Stutz, A., Melidoni, A.N., Raghunath, A., Lagreid, A., Roechert, B., Meldal, B., Aranda, B., Chen, C., Peluso, D., Galeota, E., et al. (2013). The MintAct project—IntAct as a common curation platform for 11 molecular interaction databases. *Nucleic Acids Res.* *42*, D358–D363.
- Sugiyama, N., Imamura, H., and Ishihama, Y. (2019). Large-scale discovery of substrates of the human kinome. *Sci. Rep.* *9*, 10503.
- Tanguay, P.L., Rodier, G., and Meloche, S. (2010). C-terminal domain phosphorylation of ERK3 controlled by Cdk1 and Cdc14 regulates its stability in mitosis. *Biochem. J.* *428*, 103–111.
- Terfve, C.D.A., Wilkes, E.H., Casado, P., Cutillas, P.R., and Saez-Rodriguez, J. (2015). Large-scale models of signal propagation in human cells derived from discovery phosphoproteomic data. *Nat. Commun.* *6*, 8033.
- Türei, D., Korcsmáros, T., and Saez-Rodriguez, J. (2016). OmniPath: guidelines and gateway for literature-curated signaling pathway resources. *Nat. Methods* *13*, 966–967.
- Uhlén, M., Fagerberg, L., Hallström, B.M., Lindskog, C., Oksvold, P., Mardinoglu, A., Sivertsson, Å., Kampf, C., Sjödstedt, E., Asplund, A., et al. (2015). Proteomics. Tissue-based map of the human proteome. *Science* *347*, 1260419.
- UniProt Consortium, T. (2018). UniProt: the universal protein KnowledgeBase. *Nucleic Acids Res.* *46*, 2699.
- Viger, F., and Latapy, M. (2005). Fast generation of random connected graphs with prescribed degrees. *arXiv*, arXiv:cs/0502085v1.
- Ward, J.J., McGuffin, L.J., Bryson, K., Buxton, B.F., and Jones, D.T. (2004). The DISOPRED server for the prediction of protein disorder. *Bioinformatics* *20*, 2138–2139.
- Wilkes, E.H., Terfve, C., Gribben, J.G., Saez-Rodriguez, J., and Cutillas, P.R. (2015). Empirical inference of circuitry and plasticity in a kinase signaling network. *Proc. Natl. Acad. Sci. USA* *112*, 7719–7724.
- Yu, G., and He, Q.Y. (2016). ReactomePA: an R/Bioconductor package for reactome pathway analysis and visualization. *Mol. Biosyst.* *12*, 477–479.
- Zhao, Y., Bjorbaek, C., and Moller, D.E. (1996). Regulation and interaction of pp90(rsk) isoforms with mitogen-activated protein kinases. *J. Biol. Chem.* *271*, 29773–29779.
- Zhou, F.F., Xue, Y., Chen, G.L., and Yao, X. (2004). GPS: a novel group-based phosphorylation predicting and scoring method. *Biochem. Biophys. Res. Commun.* *325*, 1443–1448.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Chemicals, Peptides, and Recombinant Proteins		
Trametinib	Selleckchem	Cat#S2673
GDC-0941	Selleckchem	Cat#S1065
Deposited Data		
504 human proteins identified as protein kinases	UniProt/Swiss-Prot Protein Knowledgebase	https://www.uniprot.org/docs/pkinfam
Phosphosite quantification across 22 kinase-inhibitory conditions in MCF7 cells	Wilkes et al., 2015	PMID: 26060313
Quantification of phosphosites across 83 breast tumor samples	Mertins et al., 2016	PMID: 27251275
Tissue RNA expression data for protein kinases	GTEx Consortium, 2013 Papatheodorou et al., 2018	https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-5214/
RNA expression data	Uhlén et al., 2015	www.proteinatlas.org
Lists of human phosphosites, kinase substrates and kinase regulatory sites	Hornbeck et al., 2015	https://www.phosphosite.org/
Phosphosite functional scores	Ochoa et al., 2020	PMID: 31819260
Kinase families	Manning et al., 2002	http://kinase.com/web/current/kinbase/
Amino acid frequencies in the human proteome	UniProt Consortium, 2018	ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/reference_proteomes/Eukaryota/UP000005640_9606.fasta.gz
Phosphorylation hot-spot data	Strumillo et al., 2019	PMID: 31036831
Features for phosphosite sign predictions	Ochoa et al., 2020	PMID: 31819260
<i>In vivo</i> kinase substrate list	Hijazi et al., 2020	PMID: 31959955
<i>In vitro</i> kinase substrate list	Sugiyama et al., 2019	PMID: 31324866
human protein-protein interaction network	Stutz et al., 2013	https://www.ebi.ac.uk/intact/
edges between PI3K and network	Kanehisa 2019; Kanehisa and Goto 2000; Kanehisa et al., 2019	https://www.genome.jp/kegg/pathway.html
Phosphoproteomic measurements for MEK1, PI3Ki and control condition	This paper	N/A
Experimental Models: Cell Lines		
Kasumi-1	ATCC	C#CRL-2724; RRID: CVCL_0589
Software and Algorithms		
e1071 v 1.6.8	N/A	https://CRAN.R-project.org/package=e1071
Omnipath Python module	Türei et al., 2016	https://github.com/saezlab/pypath
bartMachine v 1.2.3	Kapelner and Bleich 2016	https://cran.r-project.org/web/packages/bartMachine/
igraph v 1.2.2	Csárdi and Nepusz 2006	https://cran.r-project.org/web/packages/igraph/
ReactomePA v 1.22.0	Yu and He 2016	https://bioconductor.org/packages/release/bioc/html/ReactomePA.html
Reactome.db v 1.62	N/A	https://bioconductor.statistik.tu-dortmund.de/packages/3.6/data/annotation/html/reactome.db.html
ROCR v 1.0.7	Sing et al., 2005	https://cran.r-project.org/web/packages/ROCR/
Limma v 3.40.6	Ritchie et al., 2015	https://bioconductor.org/packages/release/bioc/html/limma.html
Other		
Code for predictions, figures and analysis	This paper	https://github.com/evocellnet/kinase-activity-net

LEAD CONTACT AND MATERIALS AVAILABILITY

Further information and requests for resources should be directed to and will be fulfilled by the Lead Contact, Pedro Beltrao (pbeltrao@ebi.ac.uk). This study did not generate new unique reagents.

EXPERIMENTAL MODEL AND SUBJECT DETAILS

Kasumi-1 cells (Male) were routinely cultured using RPMI (10% FBS, 1% penicillin/streptomycin) at 37°C, in a humidified atmosphere containing 5% CO₂. Each kinase inhibitor was diluted to 1000 times the desired concentration for treatment using DMSO. Cells were incubated at a density of 0.5x10⁶ cells/mL and the following day Trametinib or GDC-0941 were added at 0.5 μM for one hour prior to lysis.

METHOD DETAILS

Data

We defined the human kinome as the list of 504 human proteins identified as protein kinases in the UniProt/Swiss-Prot Protein Knowledgebase, *pkinfam* (accessed 8 November 2017 at <https://www.uniprot.org/docs/pkinfam>). Quantitative phosphoproteomic data were retrieved from two publications. The first included phosphosite quantifications of 213 phosphosites for 100 kinases across 22 kinase-inhibitory conditions in MCF7 cells (Wilkes et al., 2015). The second quantified 1537 phosphosites on 193 kinases across 83 breast tumor samples (Mertins et al., 2016). Tissue RNA expression data for protein kinases were retrieved from the GTEx project (GTEx Consortium, 2013) as provided by Expression Atlas (E-MTAB-5214, timestamp 26 April 2018) (Papatheodorou et al., 2018). We furthermore retrieved tissue RNA expression data from the Human Protein Atlas project (accessed from www.proteinatlas.org 1 December 2017) (Uhlén et al., 2015). Lists of human phosphosites, kinase substrates, and kinase regulatory sites were retrieved from the PhosphoSitePlus database (accessed May 1, 2018) (Hornbeck et al., 2015). Amino acid frequencies in the human proteome were derived from the UniProt proteome database (UniProt Consortium, 2018). *In vitro* kinase substrate list was retrieved from a publication by Sugiyama and colleagues (Sugiyama et al., 2019). *In vivo* kinase substrate list was retrieved from a publication by Hijazi and colleagues (Hijazi et al., 2020).

Protein Kinase Specificity Models

Constructing Kinase Specificity Models

We estimated kinase specificity through the construction of position-specific scoring matrices (PSSMs) from the amino acid sequences around known substrate sites (+/-7 residues), omitting autophosphorylation sites. We required at least 10 known substrates in order to build a PSSM for a given kinase, resulting in PSSMs for 140 protein kinases. In order to reduce the influence of redundant sequences on the construction of the matrices, we employed a position-based sequence-weighting method (Henikoff and Henikoff, 1994).

Given a set of $n \geq 10$ substrate amino-acid sequences, $S = \{S_1, S_2, \dots, S_i, \dots, S_{n-1}, S_n\}$, where $S_i = \{S_{i1}, S_{i2}, \dots, S_{i14}, S_{i15}\}$ and S_{ij} represents the amino acid at position j of sequence i , we give a weight to each of amino acid a at position j as follows:

$$w(a, j) = \frac{1}{c_j \sum_{i=1}^n (S_{ij} = a)}$$

where c_j is the number of unique amino acids found in position j among the substrates in S . Next, a weight is calculated for each sequence as the sum of its position-specific residue weights:

$$W(S_i) = \sum_{j=1}^{15} w(S_{ij}, j)$$

Finally, each sequence weight was normalized by the sum of all sequence weights:

$$\widehat{W}(S_i) = \frac{W(S_i)}{\sum_{k=1}^n W(S_k)}$$

A 20 × 15 PSSM can then be constructed as follows. First, we construct matrix r , such that entry r_{aj} contains the weighted count of amino acid a at position j across the sequences in S :

$$r_{aj} = n \sum_{i=1}^n V(S_{ij}, a)$$

$$V(S_{ij}, a) = \begin{cases} \widehat{W}(S_i), & \text{if } S_{ij} = a \\ 0, & \text{otherwise} \end{cases}$$

There is a non-zero probability of observing each residue at a position in the sequence; however, at small sample sizes, we are unlikely to accurately estimate low-probability occurrences. To overcome this, we added pseudocounts based on proteome-wide amino-acid frequencies in a position-specific manner (Henikoff and Henikoff, 1996). For each column j in the PSSM, we select a number of pseudocounts, B_j , to add:

$$B_j = m \times c_j$$

where m is a tune-able parameter and c_j is defined as above. Thus, empirically constrained positions (e.g. the +1 position for proline-directed kinases) will receive fewer pseudocounts, and thus lower baseline probabilities of observing other residues, than highly variable positions. We found that our results were not strongly dependent on m , so we fixed it at 1. A 20×15 matrix of pseudocounts, b , was then calculated as follows:

$$b_{aj} = B_j \times f_a$$

where f_a is the occurrence frequency of amino acid a in the proteome. This allows us to derive an empirical matrix of probabilities, p , of observing amino acid m at position j :

$$p(a,j) = \frac{b_{aj} + r_{aj}}{B_j + \sum_a r_{aj}}$$

The final PSSM was arrived at by calculating the \log_2 fold-change of p_{aj} versus the proteome-wide amino acid frequencies:

$$PWM_{aj} = \hat{p}(a,j) = \log_2 \left(\frac{p(a,j)}{f_a} \right)$$

Assigning PSSMs to Protein Kinases

In order to increase our coverage of specificity profiles to include protein kinases with few or no known substrates, we assigned to them either composite, family-wise PSSMs or PSSMs of protein kinases with similar specificity determining residues (SDRs) (Bradley et al., 2018). For each protein kinase family, we constructed a family-wise PSSM as described above using known substrates of all kinases in the family, as defined by the KinBase resource (Manning et al., 2002) (<http://kinase.com/web/current/kinbase/>). This family-wise PSSM was then assigned to any member of the family for which we could not construct a unique PSSM. PSSMs were assigned to 209 protein kinases in this manner.

Finally, for the remaining protein kinases for which no family-wise PSSM was available, we attempted to assign a PSSM based on SDR similarity. Towards this end, 10 kinase domain positions were selected, representing residues known to covary with kinase specificity and that are proximal (<4Å distance) to the kinase substrate at the active site (Bradley et al., 2018). For a given pair of kinases, sequence similarity across the 10 SDRs was calculated by summing BLOSUM62 substitution scores for each position. An 'SDR similarity' score was then calculated by dividing this sum by the maximum possible score across the 10 SDRs, such that identical kinases would yield a similarity score of 1.0.

As represented in Figure 2A, the relationship between SDR similarity and PSSM distance was explored systematically to decide upon an SDR similarity threshold to use for PSSM assignment. For this purpose, SDR similarity scores and PSSM distances were calculated for all possible pairwise comparisons of kinases with known specificity. Here, similarity between PSSMs was quantified using the Frobenius distance, which represents the sum of squared element-wise distances between matrix values, followed by taking the square-root (Ellis and Kobe, 2011). For reference, pairwise Frobenius distances were also calculated for PSSMs of the same kinase by subsampling known target sites of a given kinase, using a sample size of 25 targets sites (corresponding to the median number of target sites used for PSSM construction). The distribution of all possible pairwise distances among these 'duplicate' PSSMs had a median of 1.00 and a 97.5th percentile of 1.10 (Figure 2A, red line). We interpret PSSM distances below the 97.5th percentile to represent kinases with the same active site specificity. An SDR similarity threshold of 0.8 was therefore selected as more than half of kinase pairs above this value have PSSM distances below the 1.10 threshold. For PSSM assignment, targets from the most similar kinase(s) in the human kinome were selected, provided the SDR similarity score was above 0.8. We assigned PSSMs to a further 14 kinases through this method. For all PSSM comparisons, the phospho-acceptor column (P0: S/T/Y) was not used when calculating the Frobenius distance.

The predictive performance of family-based and SDR-based PSSM predictions was compared in Figure 2B. For every kinase of known specificity, a PSSM was assigned using the family-based and SDR-based approaches described above, and then the Frobenius distance between empirical and predicted PSSMs was calculated for both prediction methods.

Scoring Phosphosites with PSSMs

For each directed protein kinase-kinase relationship, we scored each known phosphosite on the substrate kinase using the upstream kinase's PSSM. For the +/-7 motif sequence around a given phosphosite (omitting the phosphosite itself), we calculated the PSSM score, s , as:

$$s = \sum_{j \neq 8} \hat{p}(a,j)$$

In order to make scores comparable between kinases, we then calculated a normalized score, \hat{s} , against the minimum and maximum scores attainable with the PSSM:

$$s_{min} = \sum_{j \neq 8} \hat{p}(\underset{a}{\operatorname{argmin}} \hat{p}(a, j), j)$$

$$s_{max} = \sum_{j \neq 8} \hat{p}(\underset{a}{\operatorname{argmax}} \hat{p}(a, j), j)$$

$$\hat{s} = \frac{s - s_{min}}{s_{max} - s_{min}}$$

Phosphosite Functional Scores

Predictions of functional relevance of phosphosites were retrieved from (Ochoa et al., 2020). The predictions were made on a variety of phosphosite structural, evolutionary and biochemical features. As the predictions were originally made on a strictly defined set of phosphosites derived from a reanalysis of a set of high-throughput phosphoproteomics experiments, not all of the phosphosites available in the PhosphoSitePlus database were represented. We \log_{10} -transformed the raw scores and normalized them against the minimum and maximum values to arrive at functional scores valued between 0.0 and 1.0, with larger scores reflecting a greater expectation of a functional impact of phosphorylation at that site.

Linking PSSMs to Phosphosite Functional Scores

We assessed a kinase's potential to phosphorylate a putative substrate at sites of likely functional relevance by linking the kinase's PSSM to the substrate's phosphosite functional scores via a Discounted Cumulative Gain calculation (DCG). In effect, we treat the PSSM as a "search function" and we employ the functional scores as relevance scores to determine how well a PSSM "finds" functional sites. For each substrate phosphosite with a functional score available, we calculate the PSSM score \hat{s} as above. Next, the n sites are ranked by \hat{s} in descending order, producing an associated ordering of functional scores $F = \{F_1, F_2, \dots, F_i, \dots, F_{n-1}, F_n\}$. The DCG for this kinase-substrate pair is then calculated as:

$$DCG = \sum_{i=1}^n \frac{F_i}{\log_2(i+1)}$$

Sites with higher PSSM scores, and thus lower rank i , contribute larger fractions of their functional scores to the sum. The DCG will be highest, then, if sites with high functional scores tend to have high PSSM scores.

In order to make DCG scores comparable between different kinase-substrate pairs, we normalized each score by the minimum and maximum possible DCG scores for the substrate. The minimum DCG for a substrate can be found by sorting the sites in ascending order of their functional scores; likewise, the maximum can be found by sorting the sites in descending order of their functional scores. Thus, the normalized DCG is:

$$nDCG = \frac{DCG - DCG_{min}}{DCG_{max} - DCG_{min}}$$

Coexpression and Tissue Specificity

Coexpression of protein kinases across tissues in the GTEx and Protein Atlas RNA expression datasets was calculated via Spearman's correlation after setting missing values to 0.0.

The tissue specificity of each kinase was calculated by assessing the skewness of its distribution of Protein Atlas expression values (in transcripts per million, or "TPM") across the samples, defined as

$$b_1 = \frac{m_3}{s^3} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}}$$

where x is the kinase's set of expression values across all tissues, \bar{x} is the sample mean expression value, s is the sample standard deviation, and m_3 is the third central moment of the distribution. Skewness was calculated using the "e1071" package for R (<https://CRAN.R-project.org/package=e1071>).

Phospho-Coregulation

We assessed coregulation of a pair of protein kinases by measuring the correlation between phosphorylation of their regulatory phosphosites across conditions (Wilkes et al., 2015) or tissue samples (Mertins et al., 2016) of phosphoproteomic experiments. Both experiments consisted of a table of \log_2 fold-changes for each quantified phosphosite across the conditions or samples, measuring

the relative intensities of the phosphosite as detected by mass spectrometry under each condition or sample versus a reference. The data for each experiment was quantile normalized by condition or sample (Bolstad et al., 2003).

Within an experiment, for each pair of phosphosites on two protein kinases, we calculated the correlation between fold-changes of the sites across all conditions or samples for which a quantification was available, where at least five such conditions or samples existed. We first removed any conditions in which one of the kinases was under chemical inhibition. The correlation was calculated using Spearman's rho and a p -value estimated for the correlation via the asymptotic t approximation. p -values were then $-\log_{10}$ -transformed; in the case that the estimated p was 0, we set the final value to 6. This value was then scaled by the functional scores of both sites such that only two sites with high functional scores and a high phosphorylation correlation would have a high final coregulation score. We then took the maximum such score across all site pairs as the final coregulation score for the kinase pair. Finally, the coregulation scores for all kinase pairs were normalized according to the maximum and minimum values of all pairs.

Prediction of Regulatory Relationships

Training and Validation Set

We retrieved a high-confidence set of known, directed kinase-kinase regulatory relationships from OmniPath (Türei et al., 2016) (fetched Jan 22, 2018 via the Python API). To ensure the quality of the relationships we only used those that were supported by at least two sources, providing a "positive" set of 825 relationships. It is more challenging to define a "negative" set of regulatory relationships, given the difficulty in unequivocally demonstrating a lack of regulation under all conditions. However, we assume that regulatory relationships are rare and that, given a random pair of kinases, there is unlikely to be a regulatory relationship between them. Working under this assumption, we constructed negative sets by randomly sampling from the space of possible relationships. To further reflect the presumed sparsity of the true network, we chose to construct a negative set that was 8 times larger than the positive set, which provided a slight improvement in prediction performance. This value was arbitrarily chosen to balance the diminishing performance boost from increasing negative set size with the rapidly increasing memory resources required to perform the training computations.

Feature Validation

We evaluated the performance of the following features for predicting protein kinase regulatory relationships: maximum PSSM score, maximum substrate phosphosite functional score, DCG, phosphoproteomic coregulation scores, tissue RNA coexpression, and regulator and substrate tissue expression specificity. Each feature was evaluated 100 times against a randomly sampled two-thirds of the positive set and an 8-fold larger randomized negative set.

Model Training and Prediction

In order to build a final predictive model of kinase-kinase regulatory relationships, we employed the Bayesian Additive Regression Trees (BART) method (Chipman et al., 2010). Briefly, BART is a "sum-of-trees" method, in which a series decision trees are fit to the data and used to classify data. Each tree consists of binary decision nodes reflecting a decision based on one of the features, e.g. "max. PSSM score > 0.75" or "GTEX coexpression < 0.3". The terminal nodes of the tree contain values which, once selected, contribute to the final classification value; in a sum-of-trees model, the decision values from each tree are summed to produce a value upon which this final classification is made. The BART method, in particular, uses a fixed number of trees, on which it places regularizing priors that ensure that each tree is a "weak learner", i.e. each tree contributes a small fraction of the final classification value. It does this by restricting the tree depth, shrinking terminal leaf nodes to the median, and adding noise to avoid over-fitting. Trees are fit to the data through Bayesian approaches to estimating the parameters, such as Markov-Chain Monte Carlo (MCMC) backfitting (Chipman et al., 2010).

We applied the BART model to our data as implemented in the R package "bartMachine" version 1.2.3 (Kapelner and Bleich, 2016). A notable extension of the original method provided by bartMachine is to incorporate data missingness into predictions (Kapelner and Bleich, 2015). For example, a missing phosphoproteomic coregulation value might be informative as it would indicate that phosphosites on the two kinases were never detected in the same conditions by the mass spectrometer, and thus a decision tree node asking "is the coregulation score missing?" can contribute to the final classification. We used this feature by enabling the "use_missing_data" and "use_missing_data_dummies_as_covars" parameters and disabling the "replace_missing_data_with_x_j_bar" and "impute_missingness_with_x_j_bar_for_lm" parameters. These settings in effect disable any imputation of missing data and produce new "dummy" variables that indicate whether a value is missing, which can then be incorporated in the decision trees. Model hyperparameters including the number of trees were determined using the built-in 5-fold cross-validation routine provided by the "bartMachineCV" function.

In addition to the quantitative features listed in the previous section, we also included the kinase types (serine/threonine versus tyrosine) for the regulating kinase and the substrate kinase as additional features. We evaluated the BART model on these features using the full "positive" training set and random "negative" training sets as outlined above. To this end, we performed 20 iterations of 3-fold cross-validation, using a different random "negative" set each iteration. We evaluated the true-positive rate, false-positive rate, the precision (positive predictive value) and the recall (sensitivity) of the model based on the calculated posterior probabilities assigned to the validation set. Performance metrics were calculated using the R package ROCR (Sing et al., 2005).

In order to produce our final classifications, we trained 100 different BART models to the training set, each with a different random instantiation of the "negative" set. Each model was then used to produce a posterior probability of a regulatory relationship for all kinase-kinase pairs. Finally, we took the mean of the 100 posterior probabilities for each relationship as the final classification score. For the assessment of rankings of known regulators or substrates (Figure 4C) and for the reconstruction of known pathways (Figure 5),

we applied a similar procedure. For each kinase under consideration, we built 3 different models using the “positive” set, after removing all relationships including that kinase, and random instantiations of the “negative” set. The mean posterior probability of each of that kinase’s relationships from these 3 models was used as the final prediction for the analysis.

Network Clustering and Pathway Enrichment

The resulting network was divided into clusters using the method of Blondel et al. (2008) as implemented by the R package “igraph” in the function “cluster_louvain” (Csárdi and Nepusz, 2006). This is a heuristic method that identifies clusters by optimizing modularity. The algorithm can be divided into two steps: first, a cluster is assigned to each node in the network. Next, one node i is iteratively re-assigned to each of its neighbors’ clusters and the impact on the network’s modularity is assessed. Node i is then re-assigned to the cluster where its inclusion results in the greatest gain in modularity. This process is repeated until no gain in modularity can be achieved, that is, a local maximum has been found. In the second step, a new network is constructed from the identified clusters. Edge weights between the nodes, including self-loops, are computed by summing over the weights of the links that connect nodes in each cluster. The first step is then reapplied on the resulting network. These two steps are then repeated iteratively to improve the cluster assignments.

Our aim is to predict regulatory relationships between kinases and as a result our network is directed, that is up to two directed edges connect each pair of kinases, one for each direction of regulation. As this method only clusters networks with at most one edge connecting each node pair, we retained the higher-probability edge of the two linking each pair of nodes. Prior to clustering, we removed regulatory relationships with posterior probabilities less than 0.5 in order to only retain high confidence predictions. The remaining probabilities were then max-min scaled to derive edge scores on the scale 0.0 to 1.0.

In order to determine if the derived clusters reflected known physiological relationships, we tested the clusters for enrichment in pathway annotations from the Reactome database (Fabregat et al., 2017). For the clusters with 10 or more kinases, we tested the relative frequency of pathway annotations of the kinases assigned to the cluster relative to the frequency of those annotations for the entire set of 504 kinases using the hypergeometric test as implemented by the ReactomePA package for R (Yu and He, 2016). We adjusted test p -values using the Benjamini-Hochberg method for controlling the false-discovery rate (Benjamini and Hochberg, 1995) and we set a critical value of 0.05 for testing significance. 315 kinases were annotated in Reactome V. 62 accessed through reactome.db version 1.62.0 with 6151 pathway annotations altogether.

Pathway-Annotation Distances

We extracted the human protein-protein interaction network from IntAct (version: Oct. 2018) (Stutz et al., 2013). Additionally, on this network, we integrated the human phosphorylation events extracted from SIGNOR, PhosphoSitePlus and OmniPath (Türei et al., 2016), resulting in a network containing 17089 nodes and 166757 edges. Given a pair of pathway annotations, we computed the mean of all shortest path distances between the proteins annotated for the pair.

These distances were divided into two sets: distances between pathways that are enriched in the same cluster ($n = 811$) and distances between enriched pathways across clusters ($n = 1019$). Furthermore, we excluded distances between pathways that shared kinases, which reduced our within-cluster set to 67. We used the Wilcoxon rank sum test to determine if there was a significant difference in distance between the two sets.

Network Modularity

To assess the modularity of our network we compared it to a set of randomly generated networks. Our reference network was generated by discarding all edges with probability lower than 0.5. The remaining edges were then min-max scaled to get an edge weight distribution of values between 0 and 1. A set of randomized networks ($n = 1000$) with the same degree distribution as the reference network were generated with the *sample_degseq* function in the igraph package. The “vl” method was used for network generation (Viger and Latapy, 2005). At each randomization, the edge weights of the reference network were shuffled and assigned to the randomized network. These were then clustered as described above. The modularity of the of the clustering was calculated with *modularity.igraph* as implemented in igraph (Clauset et al., 2004; Csárdi and Nepusz, 2006):

$$Q = \frac{1}{2m} \left(\sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j) \right)$$

Where m is the number of edges in the network, A is the adjacency matrix, k denotes the degree of the nodes in question and δ is an indicator function returning 0 if nodes i and j are both members of cluster c , and 1 otherwise. Applying this procedure to the random networks provided us with an empirical distribution of modularity values from which we derived an empirical p -value for the modularity of the reference network.

Prediction of Phosphosite Functional Sign

As a prerequisite for predicting the sign (activating versus inhibiting) of regulatory relationships, we first built a model to classify individual phosphosites as having either an inhibitory or activating effect on the substrate protein. As features, we used: the percentage position of the site relative to the start and end of the protein kinase domain (i.e. between 0 and 1 for sites that fall within the domain); the percentage position of the site along the protein’s length; the domain (if any) in which the phosphosite lies, including, but not

limited to, protein kinase domains; the phosphosite residue (serine/threonine or tyrosine); whether or not the substrate is a tyrosine kinase; an estimate of secondary sequence disorder, as calculated by DISOPRED (Ward et al., 2004); and the $-\log_{10}p$ -value of the site being in a phosphorylation hot-spot (Strumillo et al., 2019).

To train and validate our model, we fetched a list of human protein kinase phosphosites annotated as inducing or inhibiting activity from PhosphoSitePlus (Hornbeck et al., 2015). We built our model using BART as described above. We evaluated model performance via 20 iterations of 3-fold cross-validation on a training/validation set of 50 activating and 50 inhibiting phosphosites, sampled randomly each iteration. The final model was trained using the full training set and posterior probabilities of a phosphosite being an activating site were calculated.

We next found the probability cutoff that maximizes the Matthew's Correlation Coefficient (MCC) for classifying sites as either activating or inhibiting:

$$MCC = \frac{TA \times TI - FA \times FI}{\sqrt{(TA + FA)(TA + FI)(TI + FA)(TI + FI)}}$$

where TA , TI , FA , and FI are the numbers of true activating, true inhibiting, false activating and false inhibiting predictions at a given cutoff. Values above the cutoff were taken as "activating" predictions and those below were "inhibitory" predictions. The cutoff that maximizes the MCC was then subtracted from all predicted probabilities, yielding a score of less than zero for "inhibitory" predictions and greater than zero for "activating" predictions. Finally, these scores were rescaled so that the largest absolute value was 1 while maintaining a midpoint at zero.

Prediction of Regulatory Sign

We followed a similar procedure for classifying kinase-kinase regulatory relationships as being activating or inhibiting. The predictive features that we used were: the regulator and substrate protein kinase types (serine/threonine versus tyrosine kinases); the signed functional score; a signed formulation of the DCG; and a signed coregulation score. To derive a signed functional score, we simply assigned the sign of the phosphosite sign prediction (negative for "inhibiting", positive for "activating") to the site's functional score. We then used the signed formulation of the substrate's highest functional score as the final feature.

Signed Discounted Cumulative Gain

We modified the DCG calculation to determine whether a regulating kinase tends to "find" inhibitory or activating phosphosites on the substrate kinase. To achieve this, we applied a DCG-like calculation to the signed functional scores, where a positive sum would indicate an activating relationship and a negative sum would indicate an inhibitory one.

If the substrate phosphosites with the highest PSSM scores tend to have high functional scores with the same sign, the initial steps of the DCG cumulative sum will move in one direction. However, if the substrate has many sites and the predicted signs of the sites are unevenly distributed, the sheer number of sites alone would overcome the initial signal from the high PSSM-scoring sites. For example, if the substrate has 3 predicted inhibitory sites which all have high PSSM scores for the regulator and 10 predicted activating sites that have low PSSM scores (but high functional scores), the final DCG on the signed functional scores would ultimately be positive regardless of the site-ordering by PSSM. Thus, we formulated the signed DCG in terms of the most extreme value reached by the sum.

We begin, as with the standard DCG, by ranking the n substrate sites according to decreasing PSSM scores (\hat{s} , as described above). This produces an ordered set of their signed functional scores, $F^* = \{F_1^*, F_2^*, \dots, F_i^*, \dots, F_{n-1}^*, F_n^*\}$. We then calculate a partial DCG on this set, up to the index that produces the largest absolute sum:

$$DCG^* = \sum_{i=1}^j \frac{F_i^*}{\log_2(j+1)}$$

$$j = \underset{k}{\operatorname{argmax}} \left(\left| \sum_{i=1}^k \frac{F_i^*}{\log_2(i+1)} \right| \right) \leq n$$

We then normalize DCG^* against the most extreme value of the same sign possible for that substrate, retaining the sign. That is, if $DCG^* < 0$ we rank the substrate sites by increasing signed functional score to find DCG_{min}^* , the most extreme negative sum possible for the substrate; and otherwise we rank the sites by decreasing signed functional score to find DCG_{max}^* , the most extreme positive sum possible. Thus,

$$nDCG^* = \begin{cases} DCG^*/DCG_{max}^* & DCG^* > 0 \\ -DCG^*/DCG_{min}^* & DCG^* < 0 \end{cases}$$

Signed Coregulation Score

In order to produce a signed coregulation score, we followed the same procedure described above for the coregulation score. However, rather than using the p value of the correlation test, we used the signed correlation statistic (Spearman's rho). In order to make the test statistics comparable in spite of differing numbers of data points (i.e. number of conditions or samples in which two phosphosites have both been quantified), we z-transformed the scores:

$$z = \sqrt{\frac{n-3}{1.06}} \times \operatorname{artanh}(r)$$

where n is the number of coquantified conditions/samples for the pair of sites and r is the estimated correlation coefficient. In order to isolate correlations between likely regulatory sites, we then scaled z by the signed functional scores of the two sites. Finally, for a given pair of protein kinases, we took the most extreme scaled z value as their signed coregulation score.

In calculating these signed coregulation scores, we encountered cases that are inconsistent with a direct regulatory relationship between one protein kinase and another that is governed by a functional phosphosite on the regulator. In particular, in a direct regulatory relationship, the sign of the functional site on the regulator must be the same as the sign of the correlation. For example, if a phosphosite on the regulator is inhibitory (negative), a positive correlation of phosphorylation state with a substrate functional site could only occur through the activity of a third protein kinase (although we note that in kinases with more complicated rules of multi-site regulation, such correlations might be possible). In order to better discriminate strong signals of coregulation, we therefore removed site pairs in which the sign of the regulator's site was incoherent with the sign of the correlation.

Training and Validation of the Sign Predictor

We built a predictive model of regulatory sign from these features using BART as described above. As a training and validation set, we used 503 signed regulatory relationships (394 activating, 109 inhibitory) between protein kinases from the OmniPath database that were supported by at least two data sources. The model was validated via 20 iterations of 3-fold cross-validation, where each iteration used a different random sample of 109 activating relationships for the training/validation set.

We built 20 iterations of the final model using similar random instantiations of the training set. Finally, for each directed kinase-kinase pair, we assigned the mean posterior probability produced by these 20 models as a final regulatory sign score, where a higher value would indicate an activating relationship and a lower score would predict an inhibitory relationship. For sign prediction in the reconstruction of known pathways (Figure 5), we followed a similar "leave-one-out" procedure as described for prediction of the relationships. For each kinase under consideration, we built 3 different models after removing all relationships including that kinase from the training set. The mean posterior probability of each of that kinase's relationships being "activating" from these 3 models was used as the final prediction for the analysis.

Kinase Inhibitor Experiments

Phosphoproteomic analysis to test the predictions was carried out as described in Wilkes et al (Wilkes et al., 2015). Briefly, the Kasumi-1 cell line, growing in RPMI medium supplemented with 10% FBS, was treated with 1 μ M trametinib or GDC-0941 for 1 h. Cells were then lysed in a urea based lysis buffer. After trypsin digestion, phosphopeptides were enriched using TiO₂ chromatography and analyzed in a LS-MS/MS system consisting of an Ultimate 3000 ultra-high pressure chromatograph connected to a Q-Exactive Plus mass spectrometer. Data analysis was performed using the Mascot search engine and Pescal as described (Wilkes et al., 2015).

Identification of Down-regulated Phosphosites

By analysing phosphoproteomic data treated with *trametinib* (MEKi) and *GDC-0941* (PI3Ki), we looked for phosphosites that were down regulated by either inhibitor. We considered serine, threonine and tyrosine phosphorylated peptides even for multi-phosphorylated peptides. We log₂ transformed and quantile normalized the data to ensure that each sample followed the same distribution. To identify phosphosites that were down-regulated in each condition, we used the *limma* function as implemented by the R package *limma* (reproducibility-optimized statistical testing) (3.40.6) (Ritchie et al., 2015). Down-regulated phosphosites were selected by applying the cutoff of log₂ ratio to control of less than -1 and false discovery rate of lower than 0.1. p values were adjusted with the Benjamini-Hochberg method.

Kinase-Substrate Shortest Paths

To see if any novel pathways could be established from the perturbed kinases, we looked for the shortest path from the kinases perturbed by *trametinib* (MAPK2K1 and MAP2K2) and GDC-0941 (PI3K) to phosphosites down-regulated by their perturbation. Since PI3K is a lipid kinase we added edges between PI3K and kinases regulated by hsa:5290 (PIK3CA) and hsa:5291 (PIK3CB) and their substrate, Phosphatidylinositol-3,4,5-trisphosphate, in the KEGG database (accessed 16 October, 2019) (Kanehisa, 2019; Kanehisa and Goto, 2000; Kanehisa et al., 2019). Therefore, we added edges from PI3K to PRKCD (e.g. KEGG: hsa:04750), PRKCI (e.g. KEGG: hsa:04910), PRKCZ (e.g. KEGG: hsa:04910), SRC (e.g. KEGG: hsa:04926), AKT1 (e.g. KEGG: hsa:04151), AKT2 (e.g. KEGG: hsa:04151), ILK (e.g. KEGG: hsa:04510), MTOR (e.g. KEGG: hsa:04150/hsa04910), PDPK1 (e.g. KEGG: hsa:04150), PDPK2 (e.g. KEGG: hsa:04068), ITK (e.g. KEGG: hsa:04062) and PTK2 (e.g. KEGG: hsa:04062) were added to the network.

In order to calculate distance between perturbed kinase and down-regulated phosphosites, known kinase-substrate interactions from PhosphoSitePlus (Hornbeck et al., 2015) were added to the network as well as interactions predicted by both the *in vivo* (Hijazi et al., 2020) and *in vitro* (Sugiyama et al., 2019) experiments. Phosphosites that are known substrates of the perturbed kinases were not considered for analysis. In the case of PI3K, substrates of kinases linked to PI3K were discarded as well. We removed all edges with probability scores of less than 0.5. The function *all_shortest_paths()* as implemented in the R package *igraph* (Csárdi and Nepusz, 2006) was used to identify the shortest directed paths from the perturbed kinases to the phosphosites added to the network. The parameter *mode = "out"* was used and the edge weights were calculated by subtracting the min-max scaled edge probabilities from one. An interaction was considered novel if it was corroborated by either *in vivo* or *in vitro* experiment.

QUANTIFICATION AND STATISTICAL ANALYSES

All statistical tests and sample sizes are described in the Results section. Significance was determined at a significance level of 0.05. The tests were carried out using the R statistical computing environment (R Core Team, 2016).

DATA AND CODE AVAILABILITY

The code generated during this study is available at GitHub (<https://github.com/evocellnet/kinase-activity-net/>). The published article includes all other data generated during this study.