

Systems biology

# Co-phosphorylation networks reveal subtype-specific signaling modules in breast cancer

Marzieh Ayati <sup>1,\*</sup>, Mark R Chance<sup>2,3,4</sup> and Mehmet Koyutürk<sup>3,4,5,\*</sup>

<sup>1</sup>Department of Computer Science, University of Texas Rio Grande Valley, Edinburg, TX 78539, USA, <sup>2</sup>Department of Nutrition, <sup>3</sup>Center for Proteomics and Bioinformatics, <sup>4</sup>Case Comprehensive Cancer Center and <sup>5</sup>Department of Computer and Data Sciences, Case Western Reserve University, Cleveland, OH 44106, USA

\*To whom correspondence should be addressed.

Associate Editor: Lenore Cowen

Received on November 8, 2019; revised on July 10, 2020; editorial decision on July 19, 2020; accepted on July 22, 2020

## Abstract

**Motivation:** Protein phosphorylation is a ubiquitous mechanism of post-translational modification that plays a central role in cellular signaling. Phosphorylation is particularly important in the context of cancer, as downregulation of tumor suppressors and upregulation of oncogenes by the dysregulation of associated kinase and phosphatase networks are shown to have key roles in tumor growth and progression. Despite recent advances that enable large-scale monitoring of protein phosphorylation, these data are not fully incorporated into such computational tasks as phenotyping and subtyping of cancers.

**Results:** We develop a network-based algorithm, CoPPNet, to enable unsupervised subtyping of cancers using phosphorylation data. For this purpose, we integrate prior knowledge on evolutionary, structural and functional association of phosphosites, kinase–substrate associations and protein–protein interactions with the correlation of phosphorylation of phosphosites across different tumor samples (a.k.a co-phosphorylation) to construct a context-specific-weighted network of phosphosites. We then mine these networks to identify subnetworks with correlated phosphorylation patterns. We apply the proposed framework to two mass-spectrometry-based phosphorylation datasets for breast cancer (BC), and observe that (i) the phosphorylation pattern of the identified subnetworks are highly correlated with clinically identified subtypes, and (ii) the identified subnetworks are highly reproducible across datasets that are derived from different studies. Our results show that integration of quantitative phosphorylation data with network frameworks can provide mechanistic insights into the differences between the signaling mechanisms that drive BC subtypes. Furthermore, the reproducibility of the identified subnetworks suggests that phosphorylation can provide robust classification of disease response and markers.

**Availability and implementation:** CoPPNet is available at <http://compbio.case.edu/coppnet/>.

**Contact:** marzieh.ayati@utrgv.edu or mehmet.koyuturk@case.edu

**Supplementary information:** [Supplementary data](#) are available at *Bioinformatics* online.

## 1 Introduction

Protein phosphorylation is a ubiquitous mechanism of post-translational modification observed across cell types and species, and plays a central role in cellular signaling. Phosphorylation is regulated by networks composed of kinases, phosphatases and their substrates. Phosphorylation is particularly important in the context of cancer, as downregulation of tumor suppressors and upregulation of oncogenes (often kinases themselves) by dysregulation of the associated kinase and phosphatase networks are shown to have key roles in tumor growth and progression (Halim *et al.*, 2013; Rosell *et al.*, 2009). To this end, characterization of signaling networks enables exploration of the interconnected targets leading to the development of kinase inhibitors to treat a variety of cancers

(Butrynski *et al.*, 2010; Perrotti and Neviani, 2013). In response to the growing need for large-scale monitoring of phosphorylation, advanced mass spectrometry (MS)-based phospho-proteomics technologies have exploded. These technologies enable simultaneous identification and quantification of thousands of phosphopeptides and phosphosites from a given sample (Yates III. *et al.*, 2014).

MS-based phospho-proteomics screens create a great opportunity to discover biology that may not be observed in transcriptomic and proteomic data (Archer *et al.*, 2018). Indeed, recent research shows that, as compared to gene expression, data on post-transcriptional modifications can be more useful in subtyping cancers. As a striking example, monitoring of the specific phosphorylation pathways reveals a novel breast cancer (BC) subtype that is unique to the phospho-proteomics and cannot be captured based on

DNA mutations, mRNA-level expression, or protein expression (Mertins et al., 2016).

Although phospho-proteomics provides a critical data source to model signaling pathways, systematic methods for network analysis of phospho-proteins and phosphosites are relatively scarce. Since most of the methods designed for genomics and general proteomics are not designed to handle the complexity of phospho-proteomics, phospho-proteomic analyses are often centralized at the protein level. However, due to the *many-to-one* mapping from phosphosites to proteins (i.e. each protein may have multiple phosphorylation sites), and also multi-layer annotations (e.g. regulatory function of phosphosites and kinase-phosphosite associations), novel approaches are needed to fully leverage the richness of the data. To enable analysis of phospho-proteomic data at the level of phosphorylation sites and the relationships between these sites, we propose CoPPNet, a network-based algorithm for the analysis of phospho-proteomic data, which offers the following innovations: (i) construction of a PhosphoSite Functional Association (PSFA) network that represents the functional relationship among individual phosphosites. To create PSFA network, we incorporate known structural, evolutionary and functional associations between phosphosites, protein-protein interactions (PPIs) and kinase-substrate associations (KSA). (ii) Utilization of the PSFA network in the identification of phosphorylation modules in BC, through filtering of phosphosite pairs that are potentially functionally associated. CoPPNet accomplishes this by assigning co-phosphorylation (Co-P)-based weights to the edges in PSFA network, where Co-P quantifies the similarity of the phosphorylation patterns of phosphosites across different BC samples. We have recently introduced the notion of co-phosphorylation and used it in the context of predicting KSAs, showing that it significantly enhances the coverage and accuracy of prediction methods over those that utilize static data such as sequences, structures and generic networks (Ayati et al., 2019). Conceptually, Co-P is similar to gene co-expression, which has been shown to be effective in many biomedical applications (Liu et al., 2016; Yang et al., 2014). (iii) Development of a scoring scheme accompanied by an algorithm to identify co-phosphorylated signaling modules from this weighted PSFA network.

We test the proposed framework in the context of *unsupervised* identification of subtype-specific signaling modules in BC. For this purpose, we apply CoPPNet on two independent public phospho-proteomics datasets for BC. BC is categorized into four molecular subtypes: Luminal A, Luminal B, HER2-enriched and triple-negative (Basal-like). Among the subtypes, Luminal A has the greatest survival, and Basal has the poorest survival (Fallahpour et al., 2017). While constructing the weighted PSFA network and identifying co-phosphorylation modules on this network, we do not use any information on the samples' clinically determined subtypes.

Our results show that the statistically significant modules identified by CoPPNet are reproducible between the two independent datasets and can capture the differential phosphorylation between BC subtypes. The identified subtype-specific signaling modules have the potential to provide significant insights into the disruption of signaling processes in different cancer subtypes, and can be used in developing subtype-specific therapeutic targeting strategies for BC.

## 2 Materials and methods

The workflow of the proposed framework for unsupervised identification of co-phosphorylation (Co-P) modules is shown in Figure 1. As seen in the figure, we first construct a network to model the functional relationship between phosphorylation sites. For this purpose, we incorporate available knowledge on functional associations between phosphosites, KSAs and PPIs, and integrate these knowledge into a PSFA network. Subsequently, we use a module identification algorithm to identify subnetworks of the PSFA network that are composed of highly co-phosphorylated phosphosites (called *Co-P modules*). The premise of this approach is that, pairs of phosphosites whose phosphorylation is related to a specific cancer subtype will exhibit co-variation across different samples. For this reason, we expect that Co-P can highlight subtype-specific signaling modules even if subtype information is not available for the samples that are used to compute Co-P.

To assess the biological significance of the identified significant modules, we comprehensively evaluate their statistical significance and investigate the reproducibility of significant modules using a

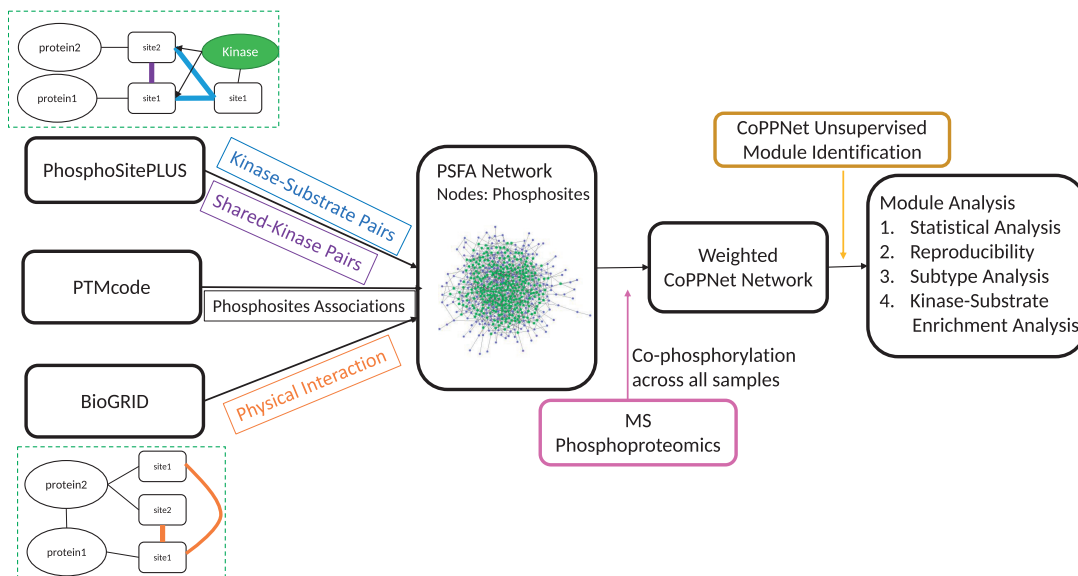


Fig. 1. Workflow of CoPPNet. We first construct a PSFA network to represent the functional relationship among phosphosites, using generic KSA, phosphosites associations and PPI data. The nodes of the PSFA network represent phosphosites and the edges represent (1) KSA, (2) phosphosites targeted by a common kinase, (3) functional associations between phosphosites, (4) physical interaction between proteins harboring the sites. For a given phosphorylation dataset collected from multiple cancer samples, we weigh the edges of the PSFA network based on the co-phosphorylation (Co-P) of pairs of sites across these samples. Then, we identify Co-P modules as subnetworks composed of heavy edges in this weighted network. Finally, we comprehensively assess the significance, reproducibility, subtype-specificity and biological relevance of the Co-P modules

dataset that comes from a different patient cohort. Subsequently, we assess the differential phosphorylation of the sites in the signaling modules between different subtypes and perform pathway enrichment analysis and kinase enrichment analysis on these modules to annotate the modules.

**PhosphoSite Functional Association (PSFA) Network.** We define a PSFA network as a network that represents *potential* functional relationships between pairs of phosphosites. This network serves the purpose of filtering out the search space for pairs of phosphosites whose co-phosphorylation may reveal their functional relationship in the context of a specific process (e.g. dysregulation of a signaling pathway in the progression of a certain cancer subtype). In PSFA network  $G(V, E)$ ,  $V$  denotes the set of nodes in the network, each of which represents a phosphosite; thus a protein is represented by multiple nodes in the PSFA network. The edge set  $E$  denotes the set of pairwise functional relationships between phosphosites, where an edge  $s_i s_j \in E$  between phosphosites  $s_i, s_j \in V$  may represent one of the following relationships:

- **Functional, Evolutionary and Structural Association between Phosphosites (FES).** PTMCode is a database of known and predicted functional associations between phosphorylation and other post-translational modification sites (Minguez et al., 2015). The associations included in PTMCode are curated from the literature, inferred from residue co-evolution, or are based on the structural distances between phosphosites. We use PTMcode as a direct source of functional, evolutionary and structural associations between phosphorylation sites.
- **Kinase-substrate association (KSA).** If phosphosite  $s_i$  is a target of kinase  $p_k$  and  $s_j$  is a phosphosite on kinase  $p_k$ , then there is an edge between  $s_i$  and  $s_j$  in the PSFA network. We call these edges *KSA edges*. This relationship indicates potential functional association between  $s_i$  and  $s_j$  since the regulation of kinase  $p_k$  through phosphorylation of  $s_j$  may influence action of  $p_k$  on  $s_i$ . In our experiments, we use PhosphositePLUS as the main source of information for KSA (Hornbeck et al., 2015).
- **Phosphosites targeted by common kinase (TCK).** If phosphosites  $s_i$  and  $s_j$  (which may be on the same protein or on different proteins) are targeted by kinase  $p_k$ , then we call them a *shared-kinase pair* and include an edge between  $s_i$  and  $s_j$  in the PSFA network. We call these edges TCK edges. We include TCK edges in the PSFA network since the activity of  $p_k$  in a specific process may influence the phosphorylation of both  $s_i$  and  $s_j$ , which may be captured by their co-phosphorylation. Indeed, studies have shown that the substrates of a protein kinase can have significant similarity in terms of their biological functions (Li et al., 2007).
- **Protein-protein interaction (PPI).** If two proteins  $p_\ell$  and  $p_r$  physically interact, for any site  $s_i$  is on  $p_\ell$ , and site  $s_j$  is on protein  $p_r$ , then there is an edge between  $s_i$  and  $s_j$  in the PSFA network. We call these edges *PPI edges*. We include PPI edges in the PSFA network, since these edges may capture functional relationships and post-transcriptional modifications beyond phosphorylation, and may remedy the sparse and incomplete nature of existing kinase-substrate annotations. In our experiments, we use the PPIs that are annotated as ‘physical’ in the BIOGRID PPI database (Chattri-Aryamontri et al., 2017) to infer the PPI edges in the PSFA network.

The PSFA network is a generic network of potential functional associations between pairs of phosphosites. In the next section, we discuss how to assign weights to the edges of the PSFA network to represent the co-phosphorylation of pairs of phosphosites in a specific context.

**Assessment of co-phosphorylation.** As with gene co-expression, correlated phosphorylation of phosphosites on proteins may be

indicative of their functional relationship in a specific biological context (Ayati et al., 2019). Based on this premise, we use context-specific phosphorylation data, obtained from MS-based phosphoproteomics assays, to assess the co-phosphorylation (Co-P) of all pairs of phosphosites that are connected in the PSFA network. In gene co-expression analysis, Pearson’s correlation and mutual information are commonly used to assess linear and non-linear relations between the expression profiles of genes (Ballouz et al., 2015; Meyer et al., 2008). Recognizing the benefits and shortcomings of each method, Song et al. (2012) developed bi-weight mid-correlation as an alternative, and showed that it outperforms mutual information in terms of capturing biologically relevant relationships between genes, while being more robust to outliers than Pearson’s correlation. Motivated by these results, we use bi-weight mid-correlation to assess the Co-P of pairs of phosphosites.

**Identification of co-phosphorylation modules.** Given a weighted PSFA network  $G(V, E, w)$  associated with a specific phosphoproteomic dataset, our objective is to identify subnetworks of the PSFA network that are enriched in highly co-phosphorylated (positively or negatively) pairs of phosphosites. This problem is similar to the well-studied problem of identifying altered subnetworks, in which the nodes are scored based on their dysregulation (e.g. z-score indicating differential gene expression) in a given condition (Dittrich et al., 2008) or association with a disease (e.g.  $-\log$  of the  $P$ -value of association) (Ayati et al., 2015). In this network, one or more connected subnetworks composed of high-scoring nodes are sought. In contrast, in our problem, scores are associated with edges, thus the problem is also similar to the infamous community detection problem in network analysis.

As with the altered subnetwork identification problem, the key component of a solution to the problem is the definition of an objective function for scoring a given subnetwork. Inspired by Newman’s definition of network modularity (Clauset et al., 2004) and our adaptation of this measure to the identification of disease-associated modules (Ayati et al., 2015), we here propose a modularity-based approach to scoring co-phosphorylation modules. In this approach, subnetworks are scored based on the difference between their total edge weight and their expected total edge weight under a reference model that considers the degree distribution of the network (in our case, the distribution of Co-P across the network). Namely, for a given set of phosphosites  $Q \subseteq V$ , we define the Co-P score of  $Q$  according to  $\ell$  as

$$\sigma(Q) = \sum_{s_i, s_j \in Q} w(s_i, s_j) - \bar{w} \quad (1)$$

where  $\bar{w}$  is the mean of the absolute values of Co-P across all pairs of phosphosites, and  $w(s_i, s_j)$  is Co-P of  $s_i$  and  $s_j$  if there is an edge between them and 0 otherwise. Thus, it penalizes the non-existent edges.

Having defined the Co-P score of a subnetwork as in Equation 1, given weighted PSFA network  $G(V, E, w)$ , we search for subnetworks of  $G$  that maximize  $\sigma(Q)$ . Since the *maximum-weight-induced subgraph problem* is NP-hard (Koyutürk et al., 2006), we use a greedy algorithm for this purpose. Namely, we search the network by starting from the phosphosite with the largest fold change, repeatedly examining the phosphosites in the neighborhood of the phosphosites so far in the subnetwork, and adding to the subnetwork the phosphosites that provide the best improvement of the subnetwork score. Once we identify a subnetwork with locally maximal Co-P score, we remove this subnetwork from  $G$  and use the greedy algorithm again to identify the next subnetwork with locally maximal Co-P score. We repeat this procedure until the entire network is exhausted, and sort all of the identified subnetworks (called Co-P modules) in decreasing order of their Co-P score. The pseudocode of the algorithm is provided in Supplementary Material. We also compare the performance of this algorithm against two other state-of-the-art module identification algorithms: Girvan and Newman’s algorithm for the identification of communities (Girvan and Newman, 2002) and the WGCNA algorithm for clustering gene co-expression networks (Langfelder and Horvath, 2008; Zhang and Horvath, 2005). We observe that the subnetworks identified by

other algorithms are less parsimonious and tend to be composed of sites that are on the same protein. CoPPNet identifies more parsimonious and statistically significant subnetworks by including a penalty term for non-extant edges in its objective function. Since the PFSA network does not include edges between sites on the same protein unless they are functionally associated, CoPPNet is able to identify signaling modules that span across multiple proteins. We report these results in detail in [Supplementary Material](#).

**Assessment of statistical significance.** To assess the statistical significance of all identified Co-P modules, we use two types of permutation tests. For this purpose, we use two null models: (i) randomize the weights of the edges of the PSFA network while preserving the topology of the network (thereby preserving the degree distribution of the phosphosites) to generate  $N$  permuted networks, and (ii) we permute the interactions while preserving the degree of phosphosites (we use  $N = 100$  in the experimental results reported in the next section). On each of the permuted networks, we identify and rank Co-P modules using the algorithm described in the previous section. We then assess the statistical significance of each module identified on the original network by comparing its score against the scores of the subnetworks that are ranked at least as high as itself on the permuted networks. We also visualize the scores of the identified modules in the context of these cumulative empirical distributions. We pick the modules that are statistically significant in terms of both null models for further analysis.

**Assessment of subtype specificity.** Although the weights of edges in the PSFA network are computed using co-phosphorylation (Co-P), which is agnostic to the subtypes of the samples, Co-P captures the co-variation of phosphorylation levels of phosphosites across different samples. Therefore, the identified modules have the potential to be associated with subtype-relevant mechanisms. Motivated by this insight, we investigate if the identified Co-P modules are composed of phosphosites that exhibit differential phosphorylation between cancer subtypes. For this purpose, we assess the differential phosphorylation of each phosphosite in a module between different subtypes. We use standard  $t$ -tests to compare the distribution of relative phosphorylation level (with respect to the common reference) in different subtypes.

**Assessment of predictive ability.** To assess the utility of identified modules in predicting subtypes, we train a support vector machine (SVM)-based classifier on one dataset using the sites in the significant modules as features and assess the performance of this classifier in predicting subtypes on the other dataset. We compare the performance of these module-based features against a full model (incorporating all sites) and a model that incorporates all sites that are significantly differentially phosphorylated ( $P < 0.05$ ) between subtypes on the training dataset.

**Assessment of reproducibility.** We assess the *reproducibility of identified Co-P modules* by investigating the overlap between significant modules identified on independent datasets. To assess the overlap between two Co-P modules that are identified in two independent datasets, we use standard hypergeometric test. We assess the *reproducibility of subtype specificity* by computing the correlation between the fold changes of sites in the modules with respect to subtypes across the two datasets. We assess the significance of this correlation empirically using a permutation test.

**Kinase substrate enrichment analysis.** Kinase substrate enrichment analysis (KSEA) seeks to identify kinases whose targets exhibit significantly altered phosphorylation levels in a given condition. KSEA scores each kinase based on the relative phosphorylation and dephosphorylation of its substrates (i.e. fold change). To assess the value added by Co-P modules, we perform kinase enrichment analysis by restricting KSEA to the substrates that are in the significant modules as opposed to all phosphosites that are identified in the study. To infer the differential activity of kinases between subtypes, we compare the score of kinases which are computed using the fold change of target phosphosites across samples in different subtypes. We identify the kinases that are predicted to have different activity by KSEA using all sites versus module-restricted sites and investigate the association of these kinases with survival using integrated gene

expression data and survival information of 1809 patients from the Gene Expression Omnibus (GEO) (Györfy et al., 2010).

**Protein expression analysis.** We also investigate if protein phosphorylation data provide information on cancer subtypes beyond what can be captured by protein expression. For this purpose, we utilize mass-spectrometry-based protein expression data that is obtained from the samples that are used to obtain the phosphoproteomic data used in our computational experiments. We utilize protein expression data in the following way: using the phosphoproteomic data, we identify phosphosites in Co-P modules that are significantly differentially expressed ( $P < 0.05$ ) between different subtypes. Subsequently, using proteomic data, we assess the differential expression of the proteins that harbor these significant phosphosites between different subtypes. If the protein that harbors the site is not identified in the protein expression data, we exclude them from the analysis. The result of this analysis is presented in [Supplementary Material](#).

## 3 Results and discussion

### 3.1 Datasets

**Phosphoproteomics data.** We use two independent public quantitative MS-based phospho-proteomics datasets obtained from BC Patient-Derived Xenografts (PDX).

- **Huang et al. data:** Huang et al. (2017) used isobaric tags for relative and absolute quantification (iTRAQ) to identify 56 874 phosphosites in 24 BC PDX models. The clinically determined subtypes for the samples in this dataset are *Basal* for 10 samples, *Luminal* for 9 samples and *HER2-enriched* for 5 samples. We remove phosphosites with missing intensity values in any sample. This results in intensity data for 15 780 phosphosites from 4539 proteins, where 13 840 serines, 2280 threonines and 67 tyrosines are phosphorylated. Protein expression data for all of these samples are also available.
- **Mertin et al. data:** The NCI Clinical Proteomic Tumor Analysis Consortium conducted an extensive MS-based phospho-proteomics of TCGA BC samples (Mertins et al., 2016). After selecting the subset of samples that have the highest coverage and filtering the phosphosites with missing intensity values in those tumors, the remaining data contained intensity values for 11 018 phosphosites mapping to 8304 phosphoproteins in 20 tumors. This dataset contains four *Basal*, nine *Luminal* and seven *HER2-enriched* samples.

**Functional, Evolutionary and Structural Association between Phosphosites (FES).** We use PTMcode, a database for functional associations of post-translational modifications within and between proteins (Minguez et al., 2015). The functional association between PTM sites have been reported based on the literature survey, co-evolution of sites, structural proximity and if PTMs at the same residue and location are within PTM highly enriched protein regions. For our analysis, we just focus on the functional associations between phosphorylation sites of different proteins.

**Kinase-substrate associations (KSAs).** We use PhosphoSitePLUS as a reference dataset for KSAs (Hornbeck et al., 2015). PhosphoSitePLUS reported 9699 KSA over 347 kinases.

**Protein-protein interaction (PPI) data.** We use a generic human PPI network downloaded from BioGRID database at <https://thebiogrid.org> (Chatr-Aryamontri et al., 2017). This network contains 194 639 interactions among 18 719 proteins.

The number of sites and edges in the final PSFA network and their types are shown in [Table 1](#). This result suggests that all different types of edges contribute to the functional relevance of the phosphosites in the modules. Although there are more PPI edges in the PSFA network, TCK edges play an important role in the identification of signaling modules, since these edges induce cliques in the



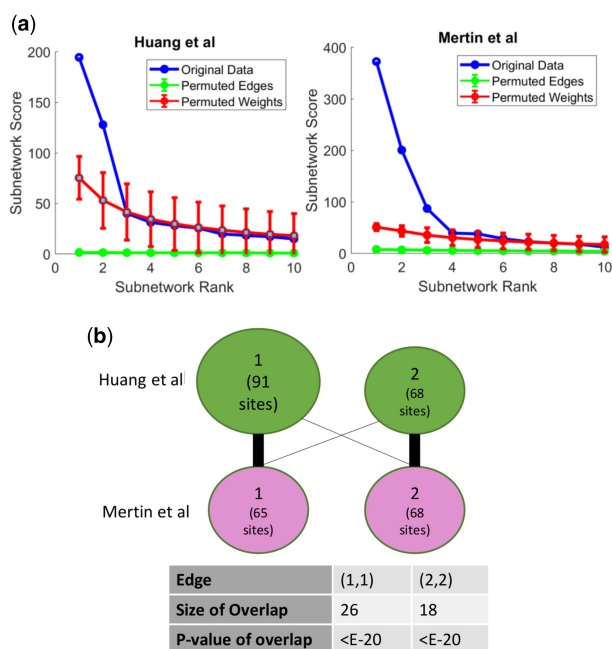
PFSA network. In this respect, CoPPNet implicitly identifies kinases whose targets exhibit enriched differential phosphorylation in specific subtypes. We elaborate on this feature of CoPPNet in the context of kinase enrichment analysis later in this section. The overlap between different types of edges is presented in [Supplementary Table S1](#).

### 3.2 CoPPNet identifies co-phosphorylation (Co-P) modules that are statistically significant and reproducible

We identify co-phosphorylated subnetworks on each of the two datasets using CoPPNet. We investigate the statistical significance of these subnetworks and visualize the results of this analysis in [Figure 2a](#). As seen in the figure, the two top-scoring subnetworks identified on both datasets have scores at least two standard deviation above the mean of the top subnetworks identified on 100 randomized networks. At a  $q$ -value threshold of 0.01, two of these subnetworks are detected to be statistically significant for each

**Table 1.** Number of phosphosites and edges in PSFA network and statistically significant modules

Type of edges # sites/# edges	PSFA network 9652/173 772	Module 191/4095	Module 268/2026
FES	7999	1	6
KSA	3024	93	306
TCK	34 857	4095	1714
PPI	133 536	46	17



**Fig. 2.** CoPPNet identifies highly significant and reproducible co-phosphorylation (Co-P) modules. (a) Statistical significance of identified subnetworks in two BC datasets. For each dataset, the blue curve shows Co-P scores (y-axis) of the highest scoring 10 subnetworks in decreasing order (rank shown on x-axis). For each rank  $i$  on the x-axis, the red (green) curve and error bar show the distribution of the scores of  $i$  highest scoring subnetworks in 100 randomized networks obtained by permutating the edge weights (edges). (b) Reproducibility of significant Co-P modules between two independent dataset Huang *et al.* and Mertin *et al.* The size of the circles indicates the number of phosphosites in each Co-P module, the number in the circle shows its rank among all identified subnetworks. The thickness of the edges represents the significance of the overlap between the two Co-P modules based on hypergeometric test. (Color version of this figure is available at [Bioinformatics](#) online.)

dataset. Note that since module identification is exhaustive, we do not expect all the identified modules to be significant. In contrast, we observe that, with the exception of the highest-scoring two modules, the scores of all other modules fall within one standard deviation of the average score of modules identified on permuted datasets. This confirms that our null models are realistic.

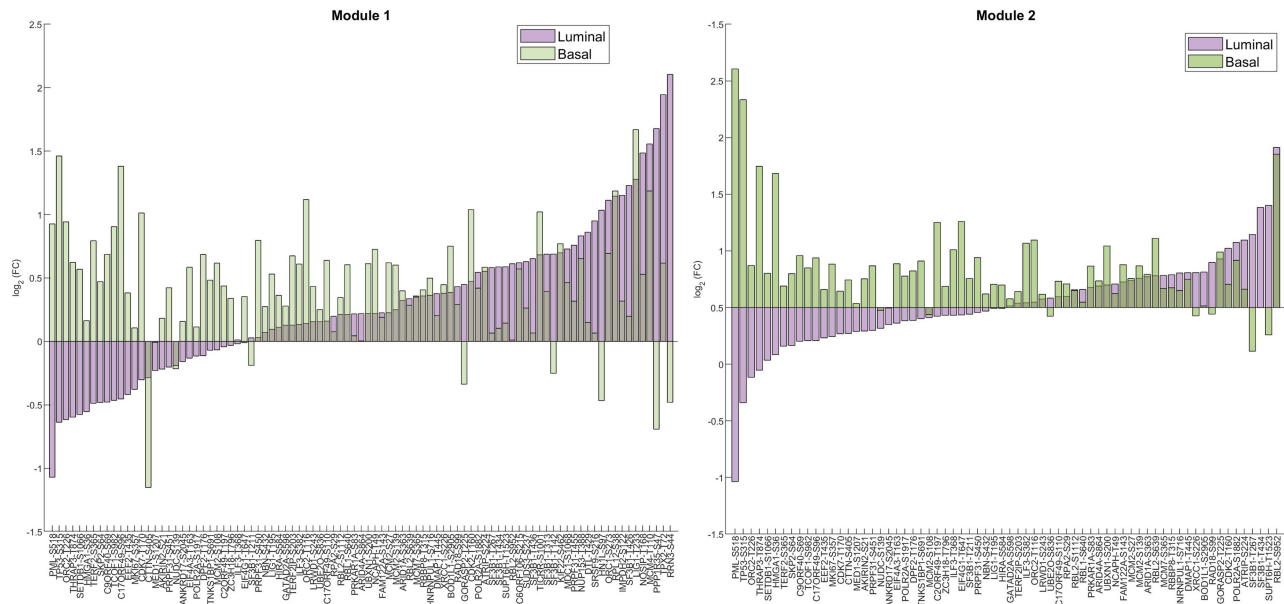
We also investigate the reproducibility of the significant modules identified on Huang *et al.* and Mertin *et al.* datasets. In [Figure 2b](#), the green circles represent the Co-P modules identified on Huang *et al.* dataset and the pink circles represent the Co-P modules identified on Mertin *et al.* dataset. As seen in the figure, there is considerable overlap between the top Co-P modules identified on each dataset; 26 out of the 91 sites in the top Huang *et al.* module and 65 sites in the top Mertin *et al.* module are identical. This overlap is highly statistically significant according to hypergeometric test and is particularly impressive considering that some phosphosites may not be present in a dataset because of the limited coverage of MS-based phospho-proteomics. Indeed, only 41 of 91 sites in the top Huang *et al.* module are identified in the Mertin *et al.* study, and only 54 of the 65 sites in the top Mertin *et al.* module are identified in the Huang *et al.* study. Many of these phospho-proteins such as THRAP3 (Beli *et al.*, 2012), NBN (Di Masi *et al.*, 2011), RAD18 (Tateishi *et al.*, 2000) and CDK7 (Li *et al.*, 2017) are playing important role in different cancers.

The second top-scoring Co-P modules identified in the two datasets, which are both highly significant ( $q < 0.01$ ), also exhibit significant overlap. Namely, 18 out of the 68 sites in the Huang *et al.* module (of which 33 are present in the Mertin *et al.* dataset) and 68 sites in the Mertin *et al.* module (of which 49 are present in the Huang *et al.* dataset) are identical. Note also that two of the sites in the top Huang *et al.* module are in the second Mertin *et al.* module, and one of the sites in the top Mertin *et al.* module is in the second-ranked Huang *et al.* module. The significant overlap and concordance between the top identified modules across two datasets show that the identified modules are highly reproducible and thus likely to be highly relevant to the dysregulation of signaling processes in BC. We also compare the significant modules with the modules extracted from gene co-expression data published in Wolf *et al.* (2014). The paper reported 11 modules. One of the co-expression modules they reported has 18 genes common with top two significant modules identified by our algorithm.

### 3.3 Co-P modules identified via unsupervised analysis are associated with BC subtypes

Since the subtype information is not used in the construction of the PSFA network and the assessment of co-phosphorylation, the identification of the Co-P modules is agnostic to the clinically determined subtypes of the samples; i.e. CoPPNet is an *unsupervised* method for the identification of BC-associated signaling modules. However, since the Co-P modules capture co-variation across different samples and this variation can be associated with subtypes, these modules can be informative on subtypes. Motivated by this consideration, we investigate if the phosphorylation levels of phosphosites in the identified modules can differentiate BC subtypes. The results of this analysis for the Huang *et al.* dataset are shown in [Figure 3](#) and [Supplementary Figure S1](#). Subtype-specific differential phosphorylation of Co-P modules identified on the Mertin *et al.* dataset are presented in [Supplementary Figure S2](#).

As seen in [Figure 3](#), top significant Co-P module identified on the Huang *et al.* dataset is highly enriched in phosphosites with significant differential phosphorylation between Luminal and Basal subtypes. There are 14 phosphosites in the top Huang *et al.* module with significant differential phosphorylation between Luminal and Basal subtypes ( $P < 0.05$ ). Eight (*DPF2-T176*, *THRAP3-T874*, *TERF2-S365*, *EIF4A3-T163*, *SETDB1-S1066*, *TCOF1-S982*, *PRPF31-S451*, *PML-S518*) out of 14 of these sites are hyper-phosphorylated in Basal samples and de-phosphorylated in Luminal samples. For some of the proteins harboring these sites, the differentiation between BC subtypes also has been captured at the level of mRNA expression. For example, *PML* (promyelocytic leukemia)



**Fig. 3.** The phosphorylation sites in top Co-P modules identified in Huang *et al.* via unsupervised analysis are associated with BC subtypes. The fold change of the phosphosites in each module are sorted in increasing order of average relative phosphorylation in Luminal samples (purple) with respect to the common reference. The green bars represent the average fold change of phosphorylation in Basal samples. (Color version of this figure is available at *Bioinformatics* online.)

**Table 2.** Performance of models for subtype prediction using all phosphosites, significant phosphosites ( $P$ -value<0.05) and phosphosites in significant Co-P modules

All sites (5476 sites) Accuracy=46%			All significant sites (621 sites) Accuracy=46%		
Real\predicted	Basal	Luminal	Real\predicted	Basal	Luminal
Basal	4	0	Basal	4	0
Luminal	7	2	Luminal	7	2

Significant sites (74 sites) Accuracy=46%			Module 1 & 2 (74 sites) Accuracy=84%		
Real\predicted	Basal	Luminal	Real\predicted	Basal	Luminal
Basal	4	0	Basal	4	0
Luminal	7	2	Luminal	2	7

and *SETDB1* (SET Domain Bifurcated 1) are significantly up-regulated in Basal cancers as compared to Luminal cancers, and their expression is related to the survival rate of the patients (Carracedo *et al.*, 2012; Jiang *et al.*, 2016). We have also compared the relative phosphorylation levels of the sites in the identified modules (Luminal versus Basal) between the Huang *et al.* and Mertin *et al.* datasets. For module 1, the Pearson, Spearman and biweight mid-correlation between the relative phosphorylation levels of the sites across the two datasets are, respectively, 0.37 ( $P < 0.004$ ), 0.37 ( $P < 0.003$ ) and 0.39 ( $P < 0.005$ ). For module 2, the Pearson, Spearman and biweight mid-correlation between the relative phosphorylation levels of the sites across the two datasets are, respectively, 0.03 ( $P < 0.41$ ), 0.41 ( $P < 0.01$ ) and 0.25 ( $P < 0.01$ ). The result of this analysis is presented in [Supplementary Figure S3](#).

### 3.4 Using co-phosphorylation modules for subtype prediction

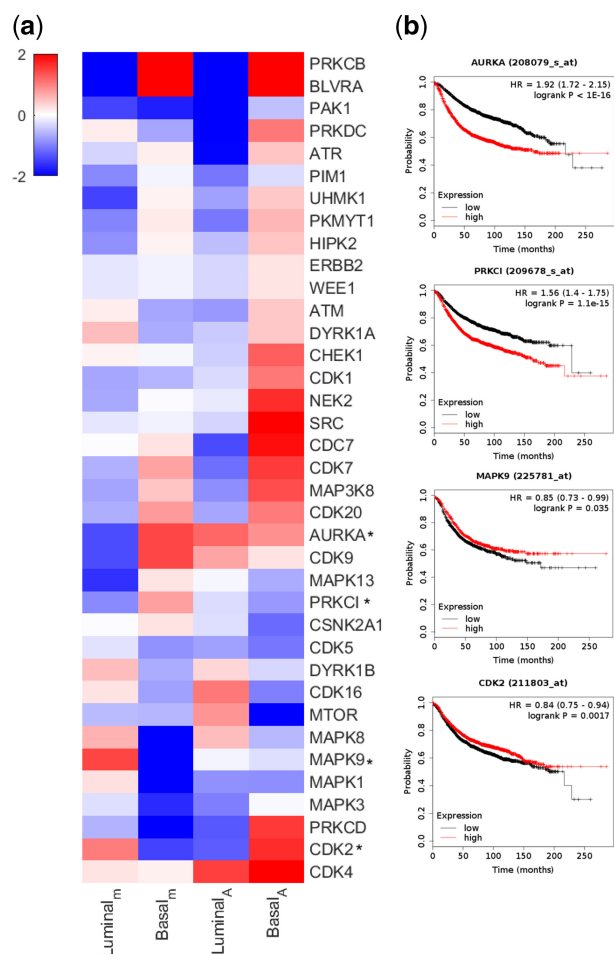
To investigate how the modules can distinguish the subtypes, we use the Co-P modules identified by CoPPNet as features for predicting subtypes on a different dataset. In this analysis, we train a SVM-based classifier for predicting subtypes using the Huang *et al.* dataset as training data. We then test the performance of this classifier on the Mertin *et al.* dataset. For this analysis, for all models that were

considered, we restricted the analysis to the sites that were identified in both datasets.

Using this setting, we compared the performance of a model which uses the sites in the significant modules (identified on training data) as features against models that use (i) all the sites that are identified in both datasets (full model), (ii) sites with significant differential phosphorylation levels ( $P < 0.05$ ) in the Huang *et al.* dataset (feature selection using significance of individual sites) and (iii) the top 74 sites according to their differential expression on the Huang *et al.* dataset (number of features identical to the number of features used by the module-based classifier). The results of this analysis are shown in [Table 2](#). As seen on the table, models that use Co-P modules outperform individual sites and significant sites. While the limited number of samples that are available pose limitations on the generalizability of these results, the improvement provided by Co-P modules demonstrates the promise of Co-P-based analysis in differentiating between subtypes.

### 3.5 Co-P modules provide a focal point for kinase activity inference

To further understand the contribution of PSFA network and co-phosphorylation analysis, we assess the value added by the Co-P modules to the inference of the differential activity of kinases



**Fig. 4.** Kinase-substrate enrichment analysis (KSEA) on Co-P modules reveals kinases that are potentially associated with BC subtype and survival. (a) The heatmap compares two different strategies for inferring kinase activity: on the left, the phosphosites used to infer kinase activity are restricted to two significant modules identified by CoPPNet (Luminal<sub>m</sub> and Basal<sub>m</sub>) on Huang *et al.* dataset. On the right, all phosphosites are used to infer kinase activity (Luminal<sub>A</sub> and Basal<sub>A</sub>). The intensity of red indicates the kinases with positive KSEA score (i.e. hyperactive in the respective subtype) and blue indicates the kinases with negative score (i.e. hypoactive in the respective subtype). Kinases that have different patterns of differential activity between subtypes in the modules versus all phosphosites are marked by a star, and their survival analysis using gene expression data is presented in (b). (Color version of this figure is available at *Bioinformatics* online.)

between Basal and Luminal subtypes. For this purpose, we use the KSEA tool, which infers the differential activity of a kinase based on the differential phosphorylation of its substrates (Casado *et al.*, 2013). In the kinase enrichment results shown in Figure 4a, the analysis is restricted to the target sites of kinases that are in the significant Co-P modules (Basal<sub>m</sub> and Luminal<sub>m</sub>) as opposed to all known target sites of the kinase that are identified in the study (Basal<sub>A</sub> and Luminal<sub>A</sub>). This analysis infers several kinases with significantly altered activity between the two subtypes. Some of these kinases show different patterns of activity when we limit the KSEA to the phosphosites in the significant modules. To assess the relevance of these kinases, we used the microarray data of BC (Györfy *et al.*, 2010), and ran Kaplan-Meier survival analysis to investigate whether the expression of these kinases is correlated with survival rate. It is well-established that Basal subtype is associated with lower survival rate as compared to Luminal subtype (Fallahpour *et al.*, 2017). We observed that, for AURKA, PRKCI, higher expression is associated with lower survival rate (Fig. 4b). KSEA analysis that is restricted to Co-P modules also suggested that these kinases are

hyperactive in the Basal samples, however, KSEA on all the phosphosites was not able to capture the association of these kinases with the subtypes. For MAPK9 and CDK2, lower expression is associated with lower survival rate, which is consistent with the kinase activity inferred by restricting to the Co-P modules. The result of this analysis for Mertin *et al.* data is presented in Supplementary Figure S5.

## 4 Conclusion

In this study, we present CoPPNet, a computational method that utilizes large-scale phospho-proteomic data for unsupervised identification of phenotype-associated signaling modules in cancer. One important contribution of the proposed method is the construction of the PSFA network which is a site-centric network that comprehensively incorporates available functional information on phosphorylation sites to enable network-based analysis of phosphorylation data. Our network model treats different types of edges identically. While observation of an edge in different databases would increase the confidence of functional association, we here use the edges only to indicate potential functional association. In future work, it can be useful to investigate the effect of assessing the value of different lines of functional evidence. Our systematic results on two BC datasets show that CoPPNet identifies reproducible subtype-specific signaling modules without requiring knowledge of the sample subtypes. However, this analysis does not account for the tissue-specificity of the phosphorylation data. Overall, this study represents one of the first attempts on utilizing phospho-proteomics to generate reproducible functional readouts of cellular signaling that can be used to characterize the dysregulation of cellular signaling in cancers and development of future therapeutic strategies.

## Acknowledgements

The authors thank Sean Maxwell, Daniella Schlatter and Ming Li for useful discussions.

## Funding

This work was supported in part by US National Institute of Health (NIH) awards R01-LM012980, R01-GM117208 and P30CA043703.

*Conflict of Interest:* none declared.

## References

- Archer, T.C. *et al.* (2018) Proteomics, post-translational modifications, and integrative analyses reveal molecular heterogeneity within medulloblastoma subgroups. *Cancer Cell*, **34**, 396–410.
- Ayati, M. *et al.* (2015) MOBAS: identification of disease-associated protein subnetworks using modularity-based scoring. *EURASIP J. Bioinf. Syst. Biol.*, **2015**, 7.
- Ayati, M. *et al.* (2019) Cophosk: a method for comprehensive kinase substrate annotation using co-phosphorylation analysis. *PLoS Comput. Biol.*, **15**, e1006678.
- Balouez, S. *et al.* (2015) Guidance for RNA-seq co-expression network construction and analysis: safety in numbers. *Bioinformatics*, **31**, 2123–2130.
- Beli, P. *et al.* (2012) Proteomic investigations reveal a role for RNA processing factor THRAP3 in the DNA damage response. *Mol. Cell*, **46**, 212–225.
- Butrynski, J.E. *et al.* (2010) Crizotinib in ALK-rearranged inflammatory myofibroblastic tumor. *N. Engl. J. Med.*, **363**, 1727–1733.
- Carracedo, A. *et al.* (2012) A metabolic prosurvival role for PML in breast cancer. *J. Clin. Investig.*, **122**, 3088–3100.
- Casado, P. *et al.* (2013) Kinase-substrate enrichment analysis provides insights into the heterogeneity of signaling pathway activation in leukemia cells. *Sci. Signal*, **6**, rs6–rs6.
- Chatr-Aryamontri, A. *et al.* (2017) The BioGRID interaction database: 2017 update. *Nucleic Acids Res.*, **45**, D369–D379.
- Clauset, A. *et al.* (2004) Finding community structure in very large networks. *Phys. Rev. E*, **70**, 066111.

- Di Masi, A. et al. (2011) Cancer predisposing mutations in BRCT domains. *IUBMB Life*, **63**, 503–512.
- Dittrich, M.T. et al. (2008) Identifying functional modules in protein–protein interaction networks: an integrated exact approach. *Bioinformatics*, **24**, i223–i231.
- Fallahpour, S. et al. (2017) Breast cancer survival by molecular subtype: a population-based analysis of cancer registry data. *CMAJ Open*, **5**, E734–E739.
- Girvan, M. and Newman, M.E. (2002) Community structure in social and biological networks. *Proc. Natl. Acad. Sci. USA*, **99**, 7821–7826.
- Györfy, B. et al. (2010) An online survival analysis tool to rapidly assess the effect of 22,277 genes on breast cancer prognosis using microarray data of 1,809 patients. *Breast Cancer Res. Treat.*, **123**, 725–731.
- Halim, V.A. et al. (2013) Comparative phosphoproteomic analysis of checkpoint recovery identifies new regulators of the DNA damage response. *Sci. Signal*, **6**, rs9.
- Hornbeck, P.V. et al. (2015) Phosphositeplus, 2014: mutations, PTMS and recalibrations. *Nucleic Acids Res.*, **43**, D512–D520.
- Huang, K.-I. et al. (2017) Proteogenomic integration reveals therapeutic targets in breast cancer xenografts. *Nat. Commun.*, **8**, 14864.
- Jiang, Y. et al. (2016) An integrated genomic analysis of Tudor domain-containing proteins identifies PHD finger protein 20-like 1 (PHF20L1) as a candidate oncogene in breast cancer. *Mol. Oncol.*, **10**, 292–302.
- Koyutürk, M. et al. (2006) Pairwise alignment of protein interaction networks. *J. Comput. Biol. J. Comput. Mol. Cell Biol.*, **13**, 182–199.
- Langfelder, P. and Horvath, S. (2008) WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics*, **9**, 559.
- Li, B. et al. (2017) Therapeutic rationale to target highly expressed CDK7 conferring poor outcomes in triple-negative breast cancer. *Cancer Res.*, **77**, 3834–3845.
- Li, T. et al. (2007) Prediction of kinase-specific phosphorylation sites with sequence features by a log-odds ratio approach. *Proteins Struct. Funct. Bioinf.*, **70**, 404–414.
- Liu, J. et al. (2016) Weighted gene co-expression network analysis identifies specific modules and hub genes related to coronary artery disease. *BMC Cardiovasc. Disord.*, **16**, 54.
- Mertins, P. et al.; NCI CPTAC. (2016) Proteogenomics connects somatic mutations to signalling in breast cancer. *Nature*, **534**, 55–62.
- Meyer, P.E. et al. (2008) minet: a R/bioconductor package for inferring large transcriptional networks using mutual information. *BMC Bioinformatics*, **9**, 461.
- Minguez, P. et al. (2015) PTMcode v2: a resource for functional associations of post-translational modifications within and between proteins. *Nucleic Acids Res.*, **43**, D494–D502.
- Perrotti, D. and Neviani, P. (2013) Protein phosphatase 2A: a target for anti-cancer therapy. *Lancet Oncol.*, **14**, e229–e238.
- Rosell, R. et al. (2009) Screening for epidermal growth factor receptor mutations in lung cancer. *N. Engl. J. Med.*, **361**, 958–967.
- Song, L. et al. (2012) Comparison of co-expression measures: mutual information, correlation, and model based indices. *BMC Bioinformatics*, **13**, 328.
- Tateishi, S. et al. (2000) Dysfunction of human rad18 results in defective post-replication repair and hypersensitivity to multiple mutagens. *Proc. Natl. Acad. Sci. USA*, **97**, 7927–7932.
- Wolf, D.M. et al. (2014) Gene co-expression modules as clinically relevant hallmarks of breast cancer diversity. *PLoS One*, **9**, e88309.
- Yang, Y. et al. (2014) Gene co-expression network analysis reveals common system-level properties of prognostic genes across cancer types. *Nat. Commun.*, **5**, 3231.
- Yates III, J. et al. (2014) Phosphoproteomics. *Analytical Anal. Chem.*, **86**, 1313–1313.
- Zhang, B. and Horvath, S. (2005) A general framework for weighted gene co-expression network analysis. *Stat. Appl. Genet. Mol. Biol.*, **4**. Article 17.