

Sequence Analysis

DeepKinZero: Zero-Shot Learning for Predicting Kinase-Phosphosite Associations Involving Understudied Kinases

Iman Deznabi^{1,2}, Busra Arabaci¹, Mehmet Koyutürk^{3, 4} and Oznur Tastan^{5,*}

¹Computer Engineering Department, Bilkent University, Ankara, 06800, Turkey

²College of Information and Computer Sciences, University of Massachusetts, Amherst, MA 01003, USA

³Dept of Computer and Data Sciences, Case Western Reserve University, Cleveland, OH 44106, USA

⁴Center for Proteomics & Bioinformatics, Case Western Reserve University, Cleveland, OH 44106, USA

⁵Faculty of Engineering and Natural Sciences, Sabanci University, Istanbul, 34956, Turkey

*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: Protein phosphorylation is a key regulator of protein function in signal transduction pathways. Kinases are the enzymes that catalyze the phosphorylation of other proteins in a target specific manner. The dysregulation of phosphorylation is associated with many diseases including cancer. Although the advances in phosphoproteomics enable the identification of phosphosites at the proteome level, most of the phosphoproteome is still in the dark: more than 95% of the reported human phosphosites have no known kinases. Determining which kinase is responsible for phosphorylating a site remains an experimental challenge. Existing computational methods require several examples of known targets of a kinase to make accurate kinase specific predictions, yet for a large body of kinases, only a few or no target sites are reported.

Results: We present DeepKinZero, the first zero-shot learning approach to predict the kinase acting on a phosphosite for kinases with no known phosphosite information. DeepKinZero transfers knowledge from kinases with many known target phosphosites to those kinases with no known sites through a zero-shot learning model. The kinase specific positional amino acid preferences are learned using a bidirectional recurrent neural network. We show that DeepKinZero achieves significant improvement in accuracy for kinases with no known phosphosites in comparison to the baseline model and other methods available. By expanding our knowledge on understudied kinases, DeepKinZero can help to chart the phosphoproteome atlas.

Availability and implementation: The source codes are available at <https://github.com/Tastanlab/DeepKinZero>.

Contact: otastan@sabanciuniv.edu

Supplementary information: Supplementary data are available at *Bioinformatics* online.

1 Introduction

Protein kinases are a large family of enzymes that catalyze the phosphorylation of other proteins (Hunter, 1995). Phosphorylation involves the transfer of a phosphoryl group to the side chain of an amino acid residue

in the substrate. The amino acid residue that receives the phosphoryl group is called the phosphorylation site, or briefly a *phosphosite*. The phosphosite is usually one of the three amino acids, serine, threonine, and tyrosine; also, other amino acids, such as histidine, are reported to act as phosphosites (Fuhs and Hunter, 2017). Phosphorylation events can lead to the activation or deactivation of proteins, modify the targets' interactions with other proteins, direct them to subcellular localization, or target them

Table 1. Kinase coverage of state-of-the-art sequence-based methods for predicting kinase-substrate associations. For each method, the middle column reports the number of kinases and kinase families for which the method can predict target phosphosites. The right column reports the criteria employed by the method for being able to make predictions for a kinase or family.

Method	Number of kinases or kinase families	Criteria for inclusion
MusiteDeep(Wang et al., 2017a)	5 families	Families with > 100 sites
PhosphoPredict(Song et al., 2017)	8 families	Families with ≥ 50 sites
Li et al.(Li et al., 2010)	8 families	Families with ≥ 50 sites
PhosphoPICK(Patrick et al., 2014)	59 human kinases	Kinases with > 10 sites
PKIS (Zou et al., 2013)	56 human kinases	Kinases with > 10 sites
KSRPred(Wang et al., 2017b)	103 human kinases	Kinases with ≥ 15 sites
KinomeExplorer (Horn et al., 2014)	222 kinases covered but accuracy assessed for 14 kinases	Kinases with ≥ 20 sites

for destruction (Pawson and Scott, 2005). Since they are the key regulators of protein function in a broad range of cellular activities, aberrant kinase function is implicated in many diseases (Gaestel et al., 2009), particularly in cancer (Blume-Jensen and Hunter, 2001; Müller et al., 2015). Several pathogenic human mutations also lie on known phosphorylation sites (Needham et al., 2019). Kinases, therefore, are also major drug targets (Klaeager et al., 2017; Ferguson and Gray, 2018). To this end, understanding the associations between kinases and phosphorylation sites holds the key to understand the signaling mechanisms in the healthy and diseased cells.

Advances in mass spectrometry-based phosphoproteomics has enabled the identification and quantification of phosphosites at the proteome level (Mann et al., 2002; Huttlin et al., 2010; Lundby et al., 2012). Many computational models have also been developed to predict phosphosites in a given input protein sequence (recently (Horn et al., 2014; Dou et al., 2014; Patrick et al., 2014; Ismail et al., 2016; Wang et al., 2017b; Song et al., 2017; Qin et al., 2016; Wang et al., 2017a) and earlier methods reviewed in (Trost and Kusalik, 2011)). Once a phosphosite is identified, either experimentally or computationally, determining the kinase that is responsible for catalyzing the phosphorylation of this site becomes the key question. With 518 identified kinases in the human genome (Manning et al., 2002) and the transient nature of kinase-substrate interactions, it is still an experimental challenge to determine the kinase that targets a given site. As underlined by a recent review (Needham et al., 2019), most of the phosphoproteome is uncharted: more than 95% of reported human phosphosites have no known kinase or associated biological function.

Several computational methods have been proposed to identify phosphorylation sites on protein sequences (Yaffe et al., 2001; Li et al., 2008; Wong et al., 2007; Qin et al., 2016; Xue et al., 2010; Koenig and Grabe, 2004; Saunders et al., 2008; Wang et al., 2017b; Song et al., 2017; Wang et al., 2017a; Blom et al., 1999; Gao et al., 2010; Patrick et al., 2014; Horn et al., 2014; Zou et al., 2013). Since these methods can also provide kinase specific predictions, they can be used to predict associated kinases of a known phosphosite. A majority of these methods utilize consensus sequence motifs or position specific scoring matrices to estimate the position preferences of each kinase. This approach requires a reasonable number of previously known targets to be able to estimate the positional preferences of a kinase accurately. Other tools employ supervised machine learning models that use a collection of established kinase-phosphosite associations. They model the relationship between the properties of kinases and the properties of their target phosphosites in a supervised classification setting. The application of such tools is limited to kinases for which a substantial number of target phosphosites are available for training. For example, MusiteDeep (Wang et al., 2017a) uses deep learning to predict binding sites for kinases, and

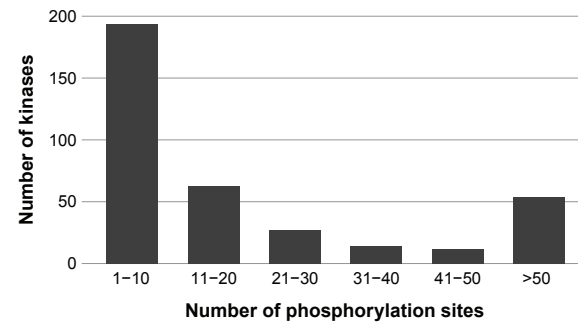


Fig. 1. The distribution of the number of experimentally validated target phosphosites for kinases in the human kinome. The histogram is based on data obtained from PhosphositePlus database experimentally validated phosphosite-kinase interactions.

it exclusively focuses on kinase families with at least 100 experimentally verified phosphosites. Recently, the use of phosphorylation data to predict kinases has been proposed, but these methods also require knowledge of target sites for a kinase to make predictions for that kinase (Ayati et al., 2019). Some of the recently developed tools and the number of kinases and/or kinase families they predict are shown in Table 1 along with the number of sites required for a kinase to be included.

A particular problem that has been overlooked in the literature is the prediction of target phosphosites for kinases with few or no known phosphosites. Despite the central role of kinases in cellular signaling cascades and their importance as potential drug targets, a large fraction of the kinome is understudied (Fedorov et al., 2010; Ferguson and Gray, 2018; Needham et al., 2019). PhosphositePlus, a database of experimentally validated phosphosites, provides phosphosite annotations for only 364 human kinases. For nearly 200 of 364 annotated kinases, there are at most 10 experimentally validated target sites (Figure 1).

In this study, we introduce DeepKinZero, a zero-shot learning approach to predict kinase-substrate associations for kinases with no known target sites. Zero-shot learning is a machine learning approach that has received significant attention, particularly in the field of computer vision. It handles recognition tasks for classes for which no training examples are available (Palatucci et al., 2009; Larochelle et al., 2008; Lampert et al., 2014; Romera-Paredes and Torr, 2015; Akata et al., 2016). The key to making predictions for classes with no training data (referred to as *unseen* or *zero-shot* classes) is to have side information which can be used to relate the classes. Based on these relations, it becomes possible to transfer the knowledge obtained from classes that have training samples (referred to as *seen* class)(Akata et al., 2016) to the previously unseen classes.

As exemplified by Yu et al. (Yu et al., 2018), it is difficult for an image classification system to recognize an okapi when there are no images of okapi in the training set. Yet, if the visual descriptions such as – zebra-stripes, four legs, brown torso, a deer-like face – can be learned from the seen classes (zebra, deer, horse, etc.) and if the system has side information indicating that okapis have these attributes, it becomes possible for the algorithm to recognize an okapi even without any prior exposure to an okapi visual. This is accomplished by detecting these visual descriptors and relating these descriptors to the side information on okapis. Similarly, even if we do not know any phosphosites that are associated with an understudied kinase (unseen class) in training, the zero-shot learning framework enables us to recognize a target site of this kinase by transferring knowledge from well-studied kinases to the rare kinases. This can be achieved by establishing a relationship between the kinases using relevant auxiliary information, such as functional, sequence, and structural characteristics of kinases. It

is important to note that, in the application of zero-shot learning to the prediction of kinase-substrate associations, phosphosites are represented as “instances” and kinases are represented as “classes” (i.e., kinase predictions are made for a given phosphosite). This is indeed the set-up that is relevant in many practical applications since the researchers who experimentally identify a phosphorylation site are interested in identifying kinases targeting that phosphosite.

Given a predicted or experimentally identified phosphosite, DeepKinZero predicts the most likely zero-shot kinase that can phosphorylate this particular site by using the local protein sequence centered at this site. DeepKinZero learns the phosphosite sequence features via a bi-directional recurrent neural network. Therein the kinases are represented based on functional and sequence information. Through learning a compatibility function that establishes relationships between the representations of the phosphosite sequences and the kinases, DeepKinZero transfers knowledge from kinases with many known phosphosites to those kinases with no known sites. We also consider alternate representations of the phosphosite sequence and the kinase embeddings and assess their effectiveness. For kinases with no known target sites (i.e., kinases for which it is not possible to make predictions using other supervised methods), DeepKinZero provides predictions with 30-fold increase in accuracy as compared to random guess.

DeepKinZero offers a scalable and flexible approach annotating sites with kinases with no prior information on their target sites. DeepKinZero is implemented in Python using Tensorflow library (Abadi et al., 2015) and is provided as an open source tool at <https://github.com/Tastanlab/DeepKinZero>.

2 Methods

2.1 Problem Formulation

The residues flanking the central phosphosite is critical for kinase specificity (Ubersax and Ferrell Jr, 2007). Thus, the local sequence surrounding the phosphorylation site has been a common input in the computational prediction of kinase-phosphosite associations. In this study, we use sequences of 15 residues (i.e., 7 residues flanking on each side of the phosphosite in addition to the phosphosite) as input and we denote these as the phosphosite sequences. Lengths of 15 or shorter have been shown to be useful in previous approaches (Hornbeck et al., 2014; Wagih et al., 2015; Trost and Kusalik, 2011). Let \mathcal{X} represent the space of phosphosite sequences and \mathcal{Y} represent the set of all identified kinases in human. The problem of kinase-phosphosite association prediction is defined as follows: given a phosphosite sequence $x \in \mathcal{X}$, identify which kinase $y \in \mathcal{Y}$ is most likely to catalyze the phosphorylation of this site. The problem is formalized as a multi-class classification problem with many classes, where each input phosphosite sequence is associated with a single kinase. This one-to-one mapping, in reality, does not always hold; a phosphosite occasionally can indeed be phosphorylated by more than one kinase. However, these cases occur rarely and in this study, whenever the predicted kinase is among the kinases known to phosphorylate a given phosphosite, we accept it as a true positive.

Some kinases are well-studied for which many target sites have been identified. On the other hand, many kinases lack formerly identified target sites. We refer to the kinases with known target sites in the training data as *common* kinases, these kinases constitute the training classes. We denote this set of kinases as $\mathcal{Y}_{tr} \subset \mathcal{Y}$. We call the kinases with few phosphosite annotation as *rare* kinases and denote the set of rare kinases as $\mathcal{Y}_{te} \subset \mathcal{Y}$. \mathcal{Y}_{te} constitutes the zero-shot test classes. By definition, the sets of common and rare kinases are disjoint, i.e., $\mathcal{Y}_{tr} \cap \mathcal{Y}_{te} = \emptyset$. **Note that the generalized zero-shot learning is a more open setting where all the classes (seen and unseen) are available as candidates for the classifier in the testing phase**

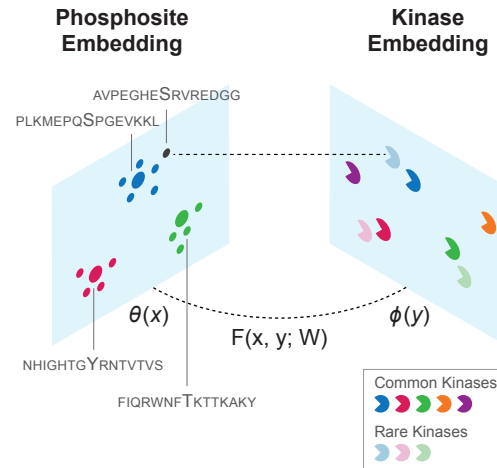


Fig. 2. Overview of the application of zero-shot learning to the prediction of kinase-phosphosite associations. The phosphosites and the kinases are embedded into multi-dimensional vector spaces using the information on sites and kinases, respectively. The training data comprise common kinases and their sites. The parameters W of the function $F(x, y; W)$ are estimated from the training data such that the compatibility between phosphosite embedding $\theta(x)$ and kinase embeddings $\phi(y)$ is maximized. For a new phosphosite at test time (shown as the black dot), the rare kinase that maximizes F for the input site’s embedding is picked by using F and the learned parameters W .

(Chao et al., 2016). This is a much harder problem which we do not tackle here and leave it as future work.

The training data contains only pairs for common kinases, $D_{tr} = \{(x_i, y_i), i = \{1, \dots, N_{tr}\}\}$, where $y_i \in \mathcal{Y}_{tr}$. Since there are no positively labeled data for the rare kinases, ($y \in \mathcal{Y}_{te}$), during the training phase, it is not possible to use traditional supervised methods to build a model for mapping sites to such rare kinases. However, it is known that some kinases are related to each other functionally, evolutionarily, or structurally (Manning et al., 2002). Thus, using zero-shot learning, the known relationships between kinases can be exploited to learn a predictive model for rare kinases. In the next section, we elaborate on this approach.

2.2 The Zero-Shot Learning Model

Following the work by Akata et al. (Akata et al., 2016), we assume that a vector space representation, called class embedding or kinase embedding, can be constructed for each kinase. Therefore, an m -dimensional “kinase embedding” vector $\phi(y) \in \mathbb{R}^m$ can be computed for each kinase $y \in \mathcal{Y}$. We expect “similar” classes to be close to each other with respect to the Euclidean metric in this embedded space. Similarly, for each phosphosite $x \in \mathcal{X}$, we compute the phosphosite embedding vector, $\theta(x) \in \mathbb{R}^d$, that represents the phosphosite sequence in a d -dimensional space. We discuss the computation of phosphosite and kinase embeddings in Sections 2.2.1 and 2.2.2 in greater detail.

The DeepKinZero Model. To accomplish transfer learning between the common and rare kinases, we learn the association between the phosphosite and the kinase embeddings. This idea is illustrated in Figure 2. Following the work in structured output prediction (Tsochantaridis et al., 2005) and prior work in zero-shot learning (Xian et al., 2017; Akata et al., 2016; Romera-Paredes and Torr, 2015; Frome et al., 2013; Akata et al., 2015; Kodirov et al., 2017; Sumbul et al., 2018), we use a compatibility function $F : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ to model the mapping between the input and output embeddings. In this model, F takes a phosphosite - kinase pair (x_i, y_j) as

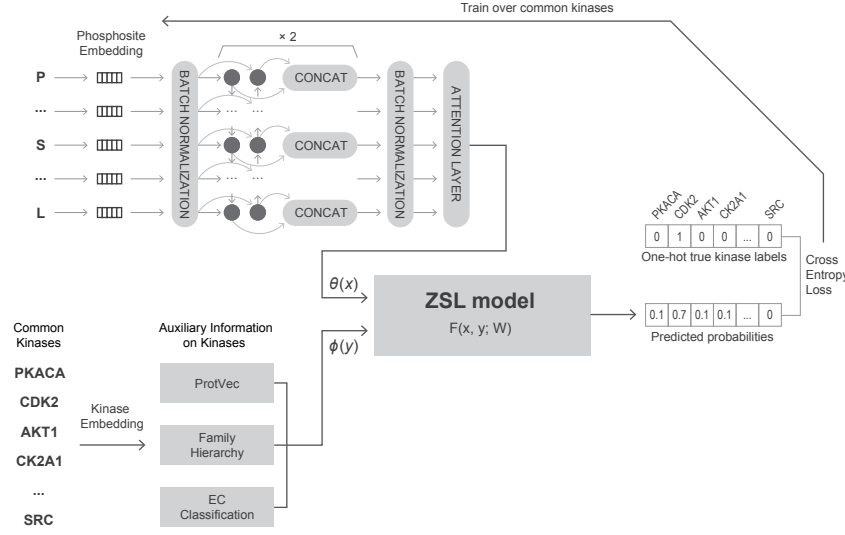


Fig. 3. The DeepKinZero framework. First, the embedded vectors of phosphosites are passed to a 2-layer bidirectional LSTM network, and then the results after passing through an attention layer are input to the ZSL model. The whole model is trained end-to-end over the common kinases.

input and returns a scalar value which is proportional to the confidence of associating the site, x_i , with kinase y_i . In this model, the probability that a given site is a target of a given kinase is calculated logarithmically from the bi-linear compatibility function F :

$$p(y|x) = \frac{\exp(F(x, y))}{\sum_{y' \in Y_{te}} \exp(F(x, y'))} \quad (1)$$

As in (Sumbul et al., 2018), we use the following bi-linear compatibility function for input x and y :

$$F(x, y) = \sum_{i=1}^d \sum_{j=1}^m W_{i,j} [\theta(x)]_i [\phi(y)]_j + \sum_{i=1}^d W_{i,m} [\theta(x)]_i + \sum_{j=1}^m W_{d,j} [\phi(y)]_j + b \quad (2)$$

which can be written in matrix notation as:

$$F(x, y) = [\theta(x)^\top \quad 1] W [\phi(y)^\top \quad 1]^\top. \quad (3)$$

Here, $[\theta(x)]_i$ and $[\phi(y)]_j$ respectively denote the i -th and the j -th entries of the phosphosite and kinase embedding vectors, respectively. W denotes the $(d+1) \times (m+1)$ compatibility matrix, where $W_{i,j}$ for $1 \leq i \leq d$ and $1 \leq j \leq m$ specifies the contribution of the correspondence between the i -th dimension in the phosphosite embedding space and the j -th dimension in the kinase embedding space to the compatibility of the phosphosite and kinase pair. $W_{d+1,i}$ and $W_{j,m+1}$ weights evaluate the information provided by the phosphosite and kinase embeddings individually. $W_{i,m+1}$ for $1 \leq i \leq d$ specifies the weight of the i -th dimension in the phosphosite embedding space, $W_{d+1,j}$ for $1 \leq j \leq m$ specifies the weight of the j -th dimension in the kinase embedding space. Finally, $W_{d+1,m+1} = b$ denotes the bias term of the model.

We represent the 15-residue phosphosite sequences centering on each phosphosite with multi-dimensional vectors in Euclidean space, such that the embeddings of similar sequences are close to each other in this space. To learn phosphosite embeddings, we use Bi-directional Recurrent Neural Network (BRNN) (Schuster and Paliwal, 1997) model with an attention mechanism over the training data. Recurrent Neural Networks

(RNNs) constitute a class of neural networks that exhibit state-of-the-art performances for modeling sequential data (Rumelhart et al., 1986). At each time step, which corresponds to the current position in the sequence, RNN accepts an input sequence vector. The hidden state of the RNN is then updated via non-linear activation functions to predict the target class, which, in our case, is the associated kinase. BRNN contains 512 LSTM cells (Hochreiter and Schmidhuber, 1997) on each direction. This number of cells is chosen to ensure the best compromise between memory requirements and accuracy performance on validation and training data.

We also employ a dot attention mechanism (Luong et al., 2015) over the output of the BRNN model to enable the model to focus on the more important positions of the input sequence. For this, we multiply the output vectors of BRNN with the attention vector A , which is $D \times 1$. D is the size of the BRNN output embeddings, which is 1024 since we have 512 nodes on each side. Let $H = [h_1, h_2, \dots, h_T]$ denote the whole output of the BRNN. To calculate the attention value for each position, we multiply the attention vector with the output vector for the position i , denoted by h_i . We apply softmax to the output of this multiplication to normalize them within the range 0-1. $\alpha_i = \text{softmax}(h_i A)$. α_i is the attention weight for position i . Finally, the phosphosite embedding vector $\phi(x)$, is the weighted average of the positions by the attention weights: $\phi(x) = \sum_{i=1}^T \alpha_i h_i$.

Training DeepKinZero. Given training data $D_{tr} = \{(x_i, y_i), i = 1, \dots, N_{tr}\}$, where $y_i \in Y_{tr}$ denote the training kinases, learning process for the zero-shot-learning model involves learning of the compatibility matrix W and the BRNN model parameters. Assuming that the training data contain independently and identically distributed samples, we estimate W that minimizes the negative log likelihood of observing the training data:

$$\hat{W} = \underset{W \in \mathbb{R}^{(d+1) \times (m+1)}}{\text{argmin}} \sum_{y_i \in Y_{tr}} -\log p(y_i | x_i) \quad (4)$$

The class posterior probabilities above are provided in Equation (1). Maximizing the likelihood is equivalent to minimizing the cross-entropy loss. We train the model end-to-end by connecting the BRNN model to ZSL model (Figure 3). In this way, the BRNN model learns phosphosite embeddings specifically useful for the ZSL model and classification of kinases. To avoid overfitting, we employ drop-out regularization with a 0.5 keep probability (Srivastava et al., 2014). We apply batch normalization in LSTM cells (Ba et al., 2016) to normalize the embeddings passed onto

the ZSL model. We initialize the W matrix randomly from a uniform distribution and minimize the cross-entropy loss function using Adam optimizer (Kingma and Ba, 2014) with learning rate 10^{-4} . The attention weights are also initialized randomly from a normal distribution with a mean of 0 and standard deviation of 0.05. The learning rate and the number of iterations are optimized on validation data (see Section 3.1 for an explanation of the validation data). To reduce the variance of the model, we ensemble 10 models each of which trained with different initializations of the model parameters. The final class probabilities are obtained by averaging output probabilities over the ensemble.

Making Predictions with DeepKinZero. The estimated \hat{W} is used at the test time; given a specific input phosphosite, the predicted kinase class, y^* , is assigned by maximizing F over the test classes:

$$y^* = \operatorname{argmax}_{y \in Y_{te}} F(x, y; \hat{W}) \quad (5)$$

This is equivalent to getting the class with the highest posterior probability as the posterior probability given in Equation (1).

2.2.1 Phosphosite Embeddings

In learning the phosphosite embeddings, we experiment inputting the phosphosite sequence with three different vector representations into the BRNN:

i) One-hot encoded vector: Each residue of a peptide sequence is coded with a 21-dimensional vector with binary entries. 20 of these dimensions encode for each of the amino acids and one extra entry is used to encode for non-extant residues. This may happen if the phosphosite is too close to the N-terminal (or the C-terminal) of the protein such that the peptide sequence is shorter than 15 residues. Eventually, with one-hot encoding, each phosphosite sequence is embedded into a $21 \times 15 = 315$ -dimensional binary vector.

ii) Physical and chemical characteristics of amino acids: We also use a reduced alphabet that represents each sequence based on the physicochemical properties of the amino acids (AA Prop) in the sequence. We consider the charge, polarity, aromaticity, size, and electronic-property of each amino acid. The categorization of each amino acid into groups based on these five properties are obtained from (Ganapathiraju et al., 2008) and is also listed in Supplementary Table 1. Using this categorization, we code each sequence based on property-based one-hot encoded vectors and concatenate them. Charge, size and aromaticity properties can each take 3 different values, polarity can take 2 and electronic property can take 5 different values. Therefore, the resulting one-hot encoded vector is $15 \times 16 = 240$ -dimensional.

iii) ProtVec: Motivated by the demonstrated success of word embedding techniques in natural language processing (e.g., Word2Vec (Mikolov et al., 2013)), unsupervised embedding models have been developed to represent protein sequences, as well. Among these models, ProtVec (Asgari and Mofrad, 2015) provides a continuous representation of protein sequences and is trained on sequences from Uniprot-SwissProt using a Skip-gram neural network (Bairoch et al., 2005). ProtVec converts each 3-gram in input sequence into a vector of length 100. There are 13 3-grams in a peptide of 15 residues, thus, our ProtVec representation of each sequence is $13 \times 100 = 1300$ -dimensional.

2.2.2 Kinase Embeddings

The key to zero-shot learning is to know, for each unseen class, the relationship with the formerly seen classes. To establish this relationship between common and rare kinases, we create four different class embedding vectors, which are then concatenated to form a kinase embedding vector, $\phi(y)$ in Equation (3). Supplementary Figure 1 summarizes the size of the kinase embedding vectors when all the sources are used. We experiment the utility of some of the vectors through computational experiments and

drop those that are not informative in the final model. Below, we give a detailed account of the sources and the way they are deployed to arrive at the desired kinase embeddings:

i) Kinase hierarchy: We use the classification proposed by (Manning et al., 2002). The data is obtained from the website Kinase.com (downloaded in June 2018). Supplementary Figure 2 shows this hierarchy. In this classification, there are 10 groups, and 116 families. We convert this to a binary vector by representing families, groups and individual kinases as one-hot encoded vectors. In the end, we attain a binary vector with a size of 583.

ii) EC classification of kinases: An alternative source of kinase categorization is the Enzyme Commission (EC) classifications provided by the ENZYME database (Bairoch, 2000) (downloaded in June 2018). According to this classification scheme, kinases are grouped into 6 main categories based on their functions. The two largest categories of kinases are the tyrosine-specific protein kinases and serine/threonine kinases. The main categories are further divided into subcategories (as shown in Supplementary Figure 3).

iii) Kin2Vec: As kinases can be related through their kinase domain sequences, we use a ProtVec representation of kinase domain sequences just as we do for the input phosphosite sequence. To differentiate the two, we refer to them as Kin2Vec. ProtVec creates vectors of length 100 for each 3-gram in the sequence and since for each kinase, the kinase domains can be of different lengths, we average the ProtVec vectors generated for each 3-gram into one vector with 100 components.

iv) KEGG pathways: To capture the relatedness of kinases in the biological functional space, we create kinase vectors based on the pathways in which the kinases participate. The human pathways are obtained from KEGG database (Kanehisa and Goto, 2000; Kanehisa et al., 2015, 2016) (downloaded in April 2018). Cumulatively, there are 190 KEGG pathways in which at least one of the kinases participate. Each kinase vector is formed as a 190-element binary vector based on its participation in each of the cellular pathways.

3 Results

3.1 Evaluation Protocol

We train and evaluate our models on the experimentally validated kinase-phosphosite associations obtained from the PhosphoSitePlus database (Hornbeck et al., 2014) (downloaded in March 2018). We exclude iso-form and fusion kinases. The dataset includes 13,426 experimentally identified phosphorylation sites and their associated 343 kinases. Following the evaluation protocol suggested by Xian et al. (Xian et al., 2017), we keep the zero-shot kinases well apart from the rest of the classes in learning the models and parameter tuning. We split the data into training, validation and test data based on the number of sites that are associated with each kinase. Kinases with more than 5 sites are considered as training classes. There are 214 such kinases. DeepKinZero is trained on this set, which contains kinase-substrate associations of 12,901 phosphorylation sites with these 214 kinases. The validation set includes the kinase-phosphosite associations of 17 kinases for which there are exactly 5 phosphorylation sites. This validation set includes 80 phosphorylation sites associated with these 17 kinases. The remaining 112 kinases with less than 5 positively labeled examples constitute the test or zero-shot classes. The test data includes these 112 kinases and kinase-phosphosite associations involving 237 phosphorylation sites.

3.2 Performance Criteria

To assess the overall performance, we use hit@k accuracy. This metric evaluates performance in terms of the number of times in which the correct

class is among the top k predicted classes, where k is a parameter. If the true class is within the top- k predicted classes, it is considered a true positive prediction. We report results for values of $k=1,3$ and 5. In cases for which a phosphosite is associated with more than one kinase, we consider the prediction to be a true positive if the model predicts one of these kinases for the corresponding phosphosite in the top k prediction. In our test dataset, 215 phosphosites are associated with a single kinase, 16 phosphosites are associated with 2 kinases, and 2 phosphosites are associated with 3 kinases. Thus, multi-class instances are rare.

3.3 Zero-Shot Learning Results

The representations of the site sequences and the kinases are critical components of the model and they can greatly influence prediction performance. For this reason, we assess the performance of DeepKinZero by comparing the prediction performance of DeepKinZero with different phosphosite and kinase embeddings.

3.3.1 The Effect of Different Phosphosite Representations on Accuracy of Predictions

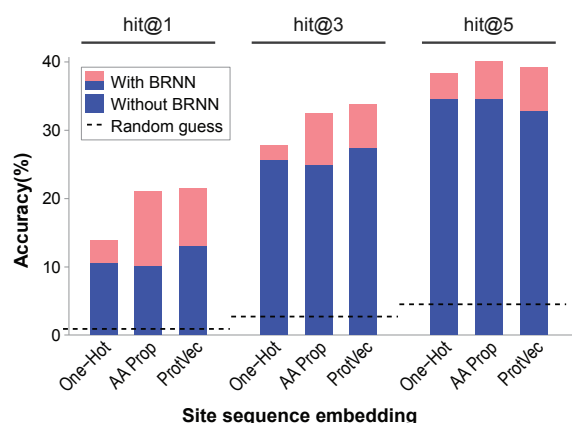


Fig. 4. The effect of phosphosite representations on the accuracy of predictions. The hit@1, hit@3, and hit@5 performance of DeepKinZero (percentage of phosphosites for which the top kinase is respectively among the top 1, 3, and 5 predictions) with six different phosphosite embedding methods (one-hot, amino-acid properties, ProtVec, each with or without a Bi-Directional Recurrent Neural Network (BRNN)) are shown. For reference, the hit@1, hit@3, and hit@5 of a random guess (the only existing alternative for the kinases tested) performance are also shown.

To thoroughly assess the effectiveness of different phosphosite representations, DeepKinZero is trained with three different input representations: One-Hot, Amino Acid Properties (AA Prop) and ProtVec with and without using BRNN. When a BRNN is employed, the BRNN is trained with the specified site sequence embeddings and the final layer of the BRNN is used as the final sequence embedding and directly input to the zero-shot classifier. Figure 4 summarizes the results of using different phosphosite sequence embeddings. As shown in Figure 4, with respect to hit@1 and hit@3 metrics, the model trained with a BRNN coupled with ProtVec vectors performs the best, where the true kinase is predicted as the top kinase for more than 20% of the sites, and it is among the top 3 for more than 30% of the sites. With respect to hit@5 metric, the input representations have less effect on the prediction performance, where amino acid properties with BRNN delivers the highest hit@5 accuracy with the true kinase being among the top 5 for more than 40% of the sites. Additionally, we observe that the use of BRNN model improves the performance. The model without BRNN embeddings that uses One-Hot sequence embedding as input only

returns the true kinase as the top prediction in 10.55% of the test cases. On the other hand, the model with BRNN and ProtVec site embeddings predict the right class with 21.52% accuracy. Note that these numbers are highly impressive since it would not be possible to train predictive models for these kinases due to the inadequacy of training samples, and random guess will achieve only 0.89% accuracy since there are 112 test classes.

To probe the usefulness of the representations learned by BRNN, we use nonlinear dimension reduction. We visualize the BRNN embeddings in a lower nonlinear dimension reduction to visualize the BRNN embeddings in a lower dimensional space using t-distributed stochastic neighbor embedding (t-SNE) (Maaten and Hinton, 2008). Supplementary Figure 4 shows that the BRNN can separate the examples in the case of kinase groups better than the ProtVec representations, hinting that it successfully captures additional critical information about kinases.

3.3.2 Effect of Kinase Embedding on Accuracy of Predictions

The performances of models trained with different kinase embeddings is shown in Table 2. In these experiments, for phosphosite embedding, we use a BRNN trained on ProtVec and compare different combinations of class embedding features with each other. To establish a baseline, the first row shows the accuracies attained using a random guess. The second row lists the performance of the model when we input the one-hot vector of kinases as class embeddings; this model is effectively a model that does not transfer knowledge between different kinases. As shown in the table, the performance of this model is worse than a random guess, demonstrating that learning is non-trivial if the class embeddings are not included. The next four rows in the table show the results of the models trained with kinase embedding vectors of individual data sources. Thus, they portray the strength of each source in isolation from the others. Among the four possible kinase embeddings, the kinase hierarchy is the leading contributor to the accuracy of the model, achieving 17.72% accuracy when used as the sole auxiliary information source. As this hierarchy reflects the functional and evolutionary information (based on sequence similarities) on the kinases, it is expected that they carry valuable information about kinase similarities. When used in isolation of other sources, Kin2Vec is found to be the least useful source.

The next set of results display the combinations of two sources. In all classes, combining family hierarchy with another information improves the model’s performance the most. The model achieves 18.99% hit@1 accuracy by combining family hierarchy with Kin2Vec. Furthermore, combining family hierarchy with EC classification or Kin2Vec vectors increases hit@5 accuracy from 37.55% to 38.82% and 40.08% respectively. Also among all combinations, its removal from the model affects the accuracy most adversely (for example second to the last row in the table).

Overall, the best performance is achieved by using family hierarchy, EC classification and Kinase2Vec vectors, which achieves 21.52% on hit@1 accuracy, 33.76% on hit@3 and 39.24% on hit@5 accuracy. Adding pathway vectors into this combination degrades the hit@1, hit@3 and hit@5 accuracies significantly, although the use of pathways alone is the second best (fourth row) when used individually as an embedding and it improves the hit@10 accuracy. It is possible that the information provided by pathway membership may not be sufficiently specific to contribute additional information on the relationships between kinases. When hit@5 or hit@10 is used, all the models except those that ignore the family hierarchy performs relatively well. The best performance is achieved when all the available information is included in the model (48.1 %).

3.3.3 Comparison with Supervised Methods Augmented with Transfer Learning

As noted before, a direct comparison between DeepKinZero and the methods which aim to predict the common kinases (Table 1) is not possible.

Table 2. **The effect of kinase embedding on the accuracy of predictions.** The hit@1, hit@3, hit@5, and hit@10 performance of DeepKinZero using all possible combinations of four different kinase embeddings are shown. Each row shows a model with a specific combination of kinase embeddings, where the check marks indicate that the corresponding kinase embedding is included in the model. For reference, the performance of random guess and an embedding that only uses the identity of individual kinases (thus does not transfer information between kinases) is also shown.

Family Hierarchy	Pathways	EC Classification	Kin2Vec	hit@1	hit@3	hit@5	hit@10
Random Guess				0.89	2.70	4.50	9.30
One-hot vector of kinases as class embedding				0.84	1.69	2.95	7.59
✓	✓	✓	✓	17.72	31.65	37.55	46.84
				8.02	13.5	16.03	21.52
				5.06	13.5	17.72	30.8
				1.27	5.91	8.02	16.03
✓	✓	✓	✓	14.77	27.85	35.02	46.84
✓				19.41	29.96	38.82	47.68
✓				18.99	33.33	40.08	47.26
	✓	✓	✓	8.86	11.39	19.41	28.69
	✓	✓	✓	8.02	13.5	16.46	21.94
		✓	✓	6.75	14.77	19.83	32.49
✓	✓	✓	✓	15.19	25.74	35.86	46.84
✓	✓	✓	✓	15.61	30.38	36.29	45.99
✓		✓	✓	21.52	33.76	39.24	47.68
	✓	✓	✓	10.55	18.14	24.05	32.91
✓	✓	✓	✓	16.88	29.11	34.18	48.10

Table 3. **Performances of augmenting existing methods with transfer learning.** Percent hit@k accuracies are given. For reference, the results achieved by DeepKinZero with the best embeddings are provided.

Model	Transfer method	hit@1	hit@3	hit@5
DeepKinZero	Zero-shot learning	21.52	33.76	39.24
PhosphoPICK(Patrick et al., 2014)	Sequence similarity	5.49	10.13	11.39
	Cosine similarity of embedding vectors	4.64	9.70	11.39
KinomeExplorer (Horn et al., 2014)	Sequence similarity	12.66	14.77	15.61
	Cosine similarity of embedding vectors	13.51	15.61	16.46

These methods will never predict the rare kinases since their candidate kinase set only comprises the common kinases. To be able to compare DeepKinZero with these methods, we develop a baseline transfer learning strategy in which we augment the traditional supervised prediction with a transfer learning step. In this baseline strategy, we first run the supervised learning method to obtain the common kinase predictions; next, we find the most similar rare kinase that shares the same family with that of the predicted kinase. We transfer the predictions within the kinase family information since this emerged as the most informative source in creating the kinase embeddings (Table 2). We finally designate this rare kinase as the method’s prediction. This comparison is only possible for methods that predict kinases as opposed to the kinase families, and we are able to apply this method to PhosphoPICK and KinomeExplorer.

To find the most similar rare kinase in the kinase family, we use two similarity assessment methods. In the first one, we pick the rare kinase that bears the highest sequence similarity to the predicted common kinase. Sequence similarity is assessed over the kinase domains global alignment (BLOSUM62, gap opening penalty of 10, and gap extension penalty of 0.5). In the second strategy, we find the closest kinase embedding vector using the cosine similarity of the kinase embedding vectors including ProtVec and EC classification vectors of the kinases (see Section 2.2.2). As can be seen, both of these results remain considerably below what DeepKinZero can achieve (see Table 3), supporting our conclusion that zero-shot learning is an effective approach to this problem.

3.3.4 Comparison with Other Phosphosite Prediction Methods for Understudied Kinases

In the literature, there are no models that we can directly compare our method against. However, there are two methods (Ellis and Kobe, 2011; Wagih et al., 2016) that aim at a different but a related problem. These two methods are designed to predict the phosphosites for kinases with no known sites, which is the reverse scenario of our problem; we predict the kinase of a given phosphosite. Predikin (Ellis and Kobe, 2011) operates with a set of rules governing the amino acids around the phosphosites. These rules, however, are derived from 3D structures of kinases bound to their substrates. Therefore, the method is limited by the availability of the protein structures and cannot be applied to kinase families without 3D structures. Because the Predikin server was not available, we were not able to carry out a comparison with this method.

The second method is by Wagih et al. (Wagih et al., 2016), which is based on the idea that, as compared to a random set of proteins, interaction partners of a kinase are more likely to be phosphorylated by that kinase. Thus, the method finds enriched motifs in the interaction partner sequences to predict sequences that a kinase can bind. The method is not applicable; however, when the kinase has a low number of interaction partners and/or the number of phosphosites on the interactors is low. Our method predicts the kinase of a given phosphosite, whereas Wagih et al. predicts the phosphosite of a kinase. Thus, the two methods are not directly comparable, but still, we conduct the following comparison. For the 112 zero-shot kinases, we predict the motifs by Wagih et al. model. If we consider the top motif returned, the method correctly matches 11 of the phosphosites of the 112 kinases, leading to 9.8% hit@1 accuracy. If we consider the top 5 motifs returned for each kinase, the correct phosphosite sequence

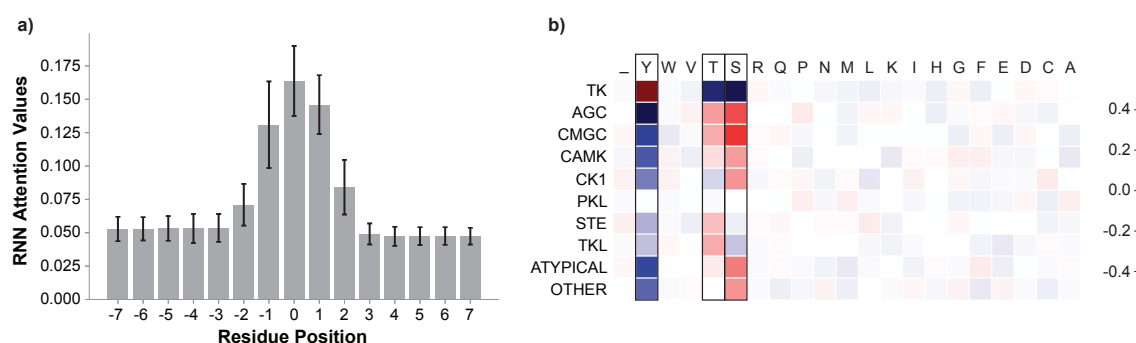


Fig. 5. Position and amino acid type weights (Best viewed in color) a) Average attention weights of the residue positions calculated over the ensemble BRNN model. Residue position 0 is the phosphosite position. b) Average zero-shot learning weights for each amino acid type at the phosphosite.

matches 26 phosphosites of the 112 kinase motifs leading to 23% hit@5 accuracy. These numbers are significantly lower than what DeepKinZero can achieve (21.52% and 40.08%). We should note that this comparison also favors Waigh et al. because DeepKinZero is evaluated based on how many phosphosites it gets right from all the available phosphosites. This is twice the number of kinases over which Waigh et al. was evaluated with.

3.3.5 Validation on an External Data

We also evaluated DeepKinZero on an external test data we had retrieved from PhosphoELM database (Diella et al., 2007) (downloaded on September 2018). We first removed all the kinases and their associated phosphosites that were in our training and validation set from PhosphositePlus dataset. The remaining kinase-substrate associations in the PhosphoELM dataset represent an instance that is well-suited to DeepKinZero’s objective, in that almost all of the kinases in this dataset have very few known associations. To be more precise, there are 52 phosphosites associated with 40 kinases and 29 of these kinases have only one site associated with it. One of them have 7 sites associated with it, while the other 10 kinases have 4, 3 or 2 associated sites. DeepKinZero trained on PhosphositePlus and evaluated on this PhosphoELM dataset achieves hit@1 accuracy of 33.96%, hit@3 accuracy of 52.83%, 62.26% hit@5 accuracy and 77.36% hit@10 accuracy. Although the dataset is small, it provides confidence that the model generalizes to other datasets.

3.4 Inspecting Model Weights

We further analyze the learned weights in the model to gain further insight into the model. First, we inspect BRNN attention weights. Figure 5 a) shows the average attention assigned to each position in the input sequence by the BRNN model. The center residue emerges as the most important residue. Thus the model correctly learns to assign more weight to the center, where the phosphosite is located. The immediate neighbors and the residues within 2 positions are the next most important residues. This aligns well with our expectations.

Next, we investigate the importance of amino acid type at the phosphosite. Recall that the W matrix specifies the relative contribution of the correspondence between each dimension in the kinase embedding space with each dimension in the site embedding space. To investigate the weights assigned to each amino acid type at the phosphosite embedding, we calculate the average weights assigned to different amino-acid types for each group of kinases at the phosphosite. As clearly seen in Figure 5 b) S, Y and T correctly receive the largest weights. Moreover, the weights assigned to different type of amino acids in each group align well with existing knowledge of kinase groups. For example, the TK family, which exclusively works on tyrosine residue (Y), puts a very large positive

weight on tyrosine while other families do not. Similarly, CMGC work predominantly on serine (S) and threonine (T) and these are the two residues that get a large positive weight. PKL group is a diverse group that could be the reason why neither of the residue types emerges as predominantly predictive.

4 Conclusion and Future Work

Many kinases are understudied with no known target proteins or sites; therefore, only a small subset of kinases dominates the annotated phosphosite databases. DeepKinZero, unlike conventional supervised methods can offer predictions for kinases which do not have any known phosphosites. The zero-shot learning framework transfers knowledge from common kinases to rare kinases, and this way, it renders the predictions for classes that were never observed in the training phase possible. Exploring the lesser-studied kinases and their associated substrates and sites will likely reveal major insights into the healthy and diseased cell. Through guiding experimental studies, we hope DeepKinZero will help in illuminating the dark phosphoproteome.

The work presented here can also be extended in different dimensions, which we plan to study in our future work. First of all, the ability to transfer learning between classes is based on the ability to define the characteristics of the kinases as vectors, which is derived from auxiliary information on kinases, such as taxonomies of kinases or deep representation of their kinase domain sequences (as detailed in Section 2.2.2 section). For a kinase to catalyze a phosphorylation event on a substrate, peptide specificity on the substrate is considered as the main determining factor. However, the peptide specificity is not the only element. The cellular localization and the structural domains outside the catalytic domain have also been reported to be important factors. Thus, in deriving the kinase embeddings, other information can be used.

We use the local peptide sequence to represent the phosphosite. Similarly, this representation can be augmented with additional structural and functional information available on the substrate. Structural features have been incorporated in kinase-substrate prediction by previous studies (Song et al., 2017), but it has been observed that these features did not significantly improve prediction performance, likely because of the limitations of training data. As more training data becomes available, transfer learning algorithms like DeepKinZero will likely enable more effective utilization of such features.

A third line of work is to extend this work to general zero-shot learning. The zero-shot learning assumes that the testing instances are only classified into the candidate unseen classes. In this study, we also assume that the candidate classes at the time of testing all belong to the rare kinases. The generalized zero-shot learning is a more open setting where all the

classes (seen and unseen) are available as candidates for the classifier at the testing phase (Chao et al., 2016). This is a much harder problem due to the greater number of classes handled during testing. Additionally, the classifier tends to assign instances into one of the previously exposed classes. This problem needs more specific methods. In future work, we plan to extend this framework to handle this generalized setting.

Acknowledgements

The authors would like to thank Dr. Ramazan Gokberk Cinbis (Middle East Technical University) for valuable discussions on zero-shot learning.

5 Funding

This work was supported by internal fundings of Sabanci and I.D. Bilkent Universities.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org.
- Akata, Z., Perronnin, F., Harchaoui, Z., and Schmid, C. (2016). Label-embedding for image classification. *IEEE transactions on pattern analysis and machine intelligence*, 38(7):1425–1438.
- Akata, Z., Reed, S., Walter, D., Lee, H., and Schiele, B. (2015). Evaluation of output embeddings for fine-grained image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2927–2936.
- Asgari, E. and Mofrad, M. R. (2015). Continuous distributed representation of biological sequences for deep proteomics and genomics. *PLoS one*, 10(11):e0141287.
- Ayati, M., Wiredja, D., Schlatter, D., Maxwell, S., Li, M., Koyutürk, M., and Chance, M. R. (2019). Cophosk: A method for comprehensive kinase substrate annotation using co-phosphorylation analysis. *PLoS Comp Biol*, 15(2).
- Ba, J. L., Kiros, J. R., and Hinton, G. E. (2016). Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Bairoch, A. (2000). The enzyme database in 2000. *Nucleic acids research*, 28(1):304–305.
- Bairoch, A., Apweiler, R., Wu, C. H., Barker, W. C., Boeckmann, B., Ferro, S., Gasteiger, E., Huang, H., Lopez, R., Magrane, M., et al. (2005). The universal protein resource (uniprot). *Nucleic acids research*, 33(suppl_1):D154–D159.
- Blom, N., Gammeltoft, S., and Brunak, S. (1999). Sequence and structure-based prediction of eukaryotic protein phosphorylation sites1. *Journal of molecular biology*, 294(5):1351–1362.
- Blume-Jensen, P. and Hunter, T. (2001). Oncogenic kinase signalling. *Nature*, 411(6835):355.
- Chao, W.-L., Changpinyo, S., Gong, B., and Sha, F. (2016). An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. In *European Conference on Computer Vision*, pages 52–68. Springer.
- Diella, F., Gould, C. M., Chica, C., Via, A., and Gibson, T. J. (2007). Phospho. elm: a database of phosphorylation sites—update 2008. *Nucleic acids research*, 36(suppl_1):D240–D244.
- Dou, Y., Yao, B., and Zhang, C. (2014). Phosphosvm: prediction of phosphorylation sites by integrating various protein sequence attributes with a support vector machine. *Amino acids*, 46(6):1459–1469.
- Ellis, J. J. and Kobe, B. (2011). Predicting protein kinase specificity: Predikin update and performance in the dream4 challenge. *PLoS one*, 6(7):e21169.
- Fedorov, O., Müller, S., and Knapp, S. (2010). The (un) targeted cancer kinome. *Nature chemical biology*, 6(3):166.
- Ferguson, F. M. and Gray, N. S. (2018). Kinase inhibitors: the road ahead. *Nature Reviews Drug Discovery*, 17(5):353.
- Frome, A., Corrado, G. S., Shlens, J., Bengio, S., Dean, J., Mikolov, T., et al. (2013). Devise: A deep visual-semantic embedding model. In *Advances in neural information processing systems*, pages 2121–2129.
- Fuhs, S. R. and Hunter, T. (2017). Phosphorylation: the emergence of histidine phosphorylation as a reversible regulatory modification. *Current opinion in cell biology*, 45:8–16.
- Gaestel, M., Kotlyarov, A., and Kracht, M. (2009). Targeting innate immunity protein kinase signalling in inflammation. *Nature Reviews Drug Discovery*, 8(6):480.
- Ganapathiraju, M., Balakrishnan, N., Reddy, R., and Klein-Seetharaman, J. (2008). Transmembrane helix prediction using amino acid property features and latent semantic analysis. In *Bmc Bioinformatics*, volume 9, page S4. BioMed Central.
- Gao, J., Thelen, J. J., Dunker, A. K., and Xu, D. (2010). Musite: a tool for global prediction of general and kinase-specific phosphorylation sites. *Molecular & Cellular Proteomics*, pages mcp–M110.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Horn, H., Schoof, E. M., Kim, J., Robin, X., Miller, M. L., Diella, F., Palma, A., Cesareni, G., Jensen, L. J., and Linding, R. (2014). Kinomexplorer: an integrated platform for kinome biology studies. *Nature methods*, 11(6):603.
- Hornbeck, P. V., Zhang, B., Murray, B., Kornhauser, J. M., Latham, V., and Skrzypek, E. (2014). Phosphositeplus, 2014: mutations, ptms and recalibrations. *Nucleic acids research*, 43(D1):D512–D520.
- Hunter, T. (1995). Protein kinases and phosphatases: the yin and yang of protein phosphorylation and signaling. *Cell*, 80(2):225–236.
- Huttlin, E. L., Jedrychowski, M. P., Elias, J. E., Goswami, T., Rad, R., Beausoleil, S. A., Villén, J., Haas, W., Sowa, M. E., and Gygi, S. P. (2010). A tissue-specific atlas of mouse protein phosphorylation and expression. *Cell*, 143(7):1174–1189.
- Ismail, H. D., Jones, A., Kim, J. H., Newman, R. H., and Kc, D. B. (2016). Rf-phos: a novel general phosphorylation site prediction tool based on random forest. *BioMed research international*, 2016.
- Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y., and Morishima, K. (2016). Kegg: new perspectives on genomes, pathways, diseases and drugs. *Nucleic acids research*, 45(D1):D353–D361.
- Kanehisa, M. and Goto, S. (2000). Kegg: kyoto encyclopedia of genes and genomes. *Nucleic acids research*, 28(1):27–30.
- Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M., and Tanabe, M. (2015). Kegg as a reference resource for gene and protein annotation. *Nucleic acids research*, 44(D1):D457–D462.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Klaeger, S., Heinzlmeir, S., Wilhelm, M., Polzer, H., Vick, B., Koenig, P.-A., Reinecke, M., Ruprecht, B., Petzoldt, S., Meng, C., et al. (2017). The target landscape of clinical kinase drugs. *Science*, 358(6367):eaan4368.
- Kodirov, E., Xiang, T., and Gong, S. (2017). Semantic autoencoder for zero-shot learning. *arXiv preprint arXiv:1704.08345*.

- Koenig, M. and Grabe, N. (2004). Highly specific prediction of phosphorylation sites in proteins. *Bioinformatics*, 20(18):3620–3627.
- Lampert, C. H., Nickisch, H., and Harmeling, S. (2014). Attribute-based classification for zero-shot visual object categorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(3):453–465.
- Larochelle, H., Erhan, D., and Bengio, Y. (2008). Zero-data learning of new tasks. In *AAAI*, volume 1, page 3.
- Li, T., Du, P., and Xu, N. (2010). Identifying human kinase-specific protein phosphorylation sites by integrating heterogeneous information from various sources. *PLoS one*, 5(11):e15411.
- Li, T., Li, F., and Zhang, X. (2008). Prediction of kinase-specific phosphorylation sites with sequence features by a log-odds ratio approach. *Proteins: Structure, Function, and Bioinformatics*, 70(2):404–414.
- Lundby, A., Secher, A., Lage, K., Nordsborg, N. B., Dmytriiev, A., Lundby, C., and Olsen, J. V. (2012). Quantitative maps of protein phosphorylation sites across 14 different rat organs and tissues. *Nature communications*, 3:876.
- Luong, M.-T., Pham, H., and Manning, C. D. (2015). Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.
- Mann, M., Ong, S.-E., Grønborg, M., Steen, H., Jensen, O. N., and Pandey, A. (2002). Analysis of protein phosphorylation using mass spectrometry: deciphering the phosphoproteome. *Trends in biotechnology*, 20(6):261–268.
- Manning, G., Whyte, D. B., Martinez, R., Hunter, T., and Sudarsanam, S. (2002). The protein kinase complement of the human genome. *Science*, 298(5600):1912–1934.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Müller, S., Chaikuad, A., Gray, N. S., and Knapp, S. (2015). The ins and outs of selective kinase inhibitor development. *Nature chemical biology*, 11(11):818.
- Needham, E. J., Parker, B. L., Burykin, T., James, D. E., and Humphrey, S. J. (2019). Illuminating the dark phosphoproteome. *Sci. Signal.*, 12(565):eaau8645.
- Palatucci, M., Pomerleau, D., Hinton, G. E., and Mitchell, T. M. (2009). Zero-shot learning with semantic output codes. In *Advances in neural information processing systems*, pages 1410–1418.
- Patrick, R., Lê Cao, K.-A., Kobe, B., and Bodén, M. (2014). Phosphopick: modelling cellular context to map kinase-substrate phosphorylation events. *Bioinformatics*, 31(3):382–389.
- Pawson, T. and Scott, J. D. (2005). Protein phosphorylation in signaling—50 years and counting. *Trends in biochemical sciences*, 30(6):286–290.
- Qin, G.-M., Li, R.-Y., and Zhao, X.-M. (2016). Phosd: inferring kinase-substrate interactions based on protein domains. *Bioinformatics*, 33(8):1197–1204.
- Romera-Paredes, B. and Torr, P. (2015). An embarrassingly simple approach to zero-shot learning. In *International Conference on Machine Learning*, pages 2152–2161.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *nature*, 323(6088):533.
- Saunders, N. F., Brinkworth, R. I., Huber, T., Kemp, B. E., and Kobe, B. (2008). Predikin and predikindb: a computational framework for the prediction of protein kinase peptide specificity and an associated database of phosphorylation sites. *BMC bioinformatics*, 9(1):245.
- Schuster, M. and Paliwal, K. K. (1997). Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11):2673–2681.
- Song, J., Wang, H., Wang, J., Leier, A., Marquez-Lago, T., Yang, B., Zhang, Z., Akutsu, T., Webb, G. I., and Daly, R. J. (2017). Phosphopredict: A bioinformatics tool for prediction of human kinase-specific phosphorylation substrates and sites by integrating heterogeneous feature selection. *Scientific Reports*, 7(1):6862.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958.
- Sumbul, G., Cinbis, R. G., and Aksoy, S. (2018). Fine-grained object recognition and zero-shot learning in remote sensing imagery. *IEEE Transactions on Geoscience and Remote Sensing*, 56(2):770–779.
- Trost, B. and Kusalik, A. (2011). Computational prediction of eukaryotic phosphorylation sites. *Bioinformatics*, 27(21):2927–2935.
- Tsochantaridis, I., Joachims, T., Hofmann, T., and Altun, Y. (2005). Large margin methods for structured and interdependent output variables. *Journal of machine learning research*, 6(Sep):1453–1484.
- Ubersax, J. A. and Ferrell Jr, J. E. (2007). Mechanisms of specificity in protein phosphorylation. *Nature reviews Molecular cell biology*, 8(7):530.
- Wagih, O., Reimand, J., and Bader, G. D. (2015). Mimp: predicting the impact of mutations on kinase-substrate phosphorylation. *Nature methods*, 12(6):531.
- Wagih, O., Sugiyama, N., Ishihama, Y., and Beltrao, P. (2016). Uncovering phosphorylation-based specificities through functional interaction networks. *Molecular & Cellular Proteomics*, 15(1):236–245.
- Wang, D., Zeng, S., Xu, C., Qiu, W., Liang, Y., Joshi, T., and Xu, D. (2017a). Musitedeep: a deep-learning framework for general and kinase-specific phosphorylation site prediction. *Bioinformatics*, 33(24):3909–3916.
- Wang, M., Wang, T., Wang, B., Liu, Y., and Li, A. (2017b). A novel phosphorylation site-kinase network-based method for the accurate prediction of kinase-substrate relationships. *BioMed research international*, 2017.
- Wong, Y.-H., Lee, T.-Y., Liang, H.-K., Huang, C.-M., Wang, T.-Y., Yang, Y.-H., Chu, C.-H., Huang, H.-D., Ko, M.-T., and Hwang, J.-K. (2007). Kinasephos 2.0: a web server for identifying protein kinase-specific phosphorylation sites based on sequences and coupling patterns. *Nucleic acids research*, 35(suppl_2):W588–W594.
- Xian, Y., Schiele, B., and Akata, Z. (2017). Zero-shot learning—the good, the bad and the ugly. *arXiv preprint arXiv:1703.04394*.
- Xue, Y., Liu, Z., Cao, J., Ma, Q., Gao, X., Wang, Q., Jin, C., Zhou, Y., Wen, L., and Ren, J. (2010). Gps 2.1: enhanced prediction of kinase-specific phosphorylation sites with an algorithm of motif length selection. *Protein Engineering, Design & Selection*, 24(3):255–260.
- Yaffe, M. B., Lepar, G. G., Lai, J., Obata, T., Volinia, S., and Cantley, L. C. (2001). A motif-based profile scanning approach for genome-wide prediction of signaling pathways. *Nature biotechnology*, 19(4):348.
- Yu, Y., Ji, Z., Guo, J., and Zhang, Z. (2018). Zero-shot learning via latent space encoding. *IEEE transactions on cybernetics*, (99):1–12.
- Zou, L., Wang, M., Shen, Y., Liao, J., Li, A., and Wang, M. (2013). Pkis: computational identification of protein kinases for experimentally discovered protein phosphorylation sites. *BMC bioinformatics*, 14(1):247.