

## RESEARCH ARTICLE

## DeepHE: Accurately predicting human essential genes based on deep learning

Xue Zhang<sup>1,2</sup>, Wangxin Xiao<sup>3\*</sup>, Weijia Xiao<sup>4</sup>

**1** Faculty of Computer and Software Engineering, Huaiyin Institute of Technology, Huai'an, Jiangsu, China, **2** School of Medicine, Tufts University, Boston, Massachusetts, United States of America, **3** Faculty of Transportation Engineering, Huaiyin Institute of Technology, Huai'an, Jiangsu, China, **4** Boston Latin School, Boston, Massachusetts, United States of America

\* [xiaozhangdoctor@126.com](mailto:xiaozhangdoctor@126.com)

## Abstract

Accurately predicting essential genes using computational methods can greatly reduce the effort in finding them via wet experiments at both time and resource scales, and further accelerate the process of drug discovery. Several computational methods have been proposed for predicting essential genes in model organisms by integrating multiple biological data sources either via centrality measures or machine learning based methods. However, the methods aiming to predict human essential genes are still limited and the performance still need improve. In addition, most of the machine learning based essential gene prediction methods are lack of skills to handle the imbalanced learning issue inherent in the essential gene prediction problem, which might be one factor affecting their performance. We propose a deep learning based method, DeepHE, to predict human essential genes by integrating features derived from sequence data and protein-protein interaction (PPI) network. A deep learning based network embedding method is utilized to automatically learn features from PPI network. In addition, 89 sequence features were derived from DNA sequence and protein sequence for each gene. These two types of features are integrated to train a multilayer neural network. A cost-sensitive technique is used to address the imbalanced learning problem when training the deep neural network. The experimental results for predicting human essential genes show that our proposed method, DeepHE, can accurately predict human gene essentiality with an average performance of AUC higher than 94%, the area under precision-recall curve (AP) higher than 90%, and the accuracy higher than 90%. We also compare DeepHE with several widely used traditional machine learning models (SVM, Naïve Bayes, Random Forest, and Adaboost) using the same features and utilizing the same cost-sensitive technique to against the imbalanced learning issue. The experimental results show that DeepHE significantly outperforms the compared machine learning models. We have demonstrated that human essential genes can be accurately predicted by designing effective machine learning algorithm and integrating representative features captured from available biological data. The proposed deep learning framework is effective for such task.

## OPEN ACCESS

**Citation:** Zhang X, Xiao W, Xiao W (2020) DeepHE: Accurately predicting human essential genes based on deep learning. *PLoS Comput Biol* 16(9): e1008229. <https://doi.org/10.1371/journal.pcbi.1008229>

**Editor:** Sonika Tyagi, Monash University, AUSTRALIA

**Received:** March 2, 2020

**Accepted:** August 9, 2020

**Published:** September 16, 2020

**Copyright:** © 2020 Zhang et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** All data used in this study are third party and freely accessible from public databases. Protein-protein interaction data are available from BioGRID database at <http://thebiogrid.org/download.php>. Essential genes data and the corresponding sequence data from DEG database are available at <http://tubic.tju.edu.cn/deg/>. DNA sequence and protein sequence data are available at [https://useast.ensembl.org/Homo\\_sapiens/Info/Annotation](https://useast.ensembl.org/Homo_sapiens/Info/Annotation). The python code is freely available at <https://github.com/xzhang2016/DeepHE>.

**Funding:** This work was supported by the National Natural Science Foundation of China, No. 61402423, XZ; National Natural Science Foundation of China, No.51678282, WXX; National Natural Science Foundation of China, No.51378243, WXX; Guizhou Provincial Science and Technology Fund with grant No. [2015]2135, XZ. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

## Author summary

Essential genes are a subset of genes. They are indispensable to the survival or reproduction of a living organism and thus play a very important role in maintaining cellular life. The identification of gene essentiality is very important for understanding the minimal requirements of an organism, identifying disease genes, and finding new drug targets. Essential genes can be identified via several wet-lab experimental methods, but these methods are often time-consuming, laborious, and costly. As a complement to the experimental methods, some centrality measures and traditional machine learning based computational methods have been proposed which mainly focused on predicting essential genes on model organisms. Here, we show that human essential genes can be accurately predicted by exploring sequence data and protein interaction network based on deep learning techniques. The ability to accurately and efficiently predict essential genes by utilizing existing biological omics data accelerates the annotation and analysis of essential genes, advance our understanding of the mechanism of basic life, and boosts the drug development.

This is a *PLOS Computational Biology Methods* paper.

## Introduction

Essential genes are a subset of genes which are indispensable to the survival or reproduction of a living organism. The prediction of gene essentiality is very important for understanding the minimal requirements of an organism, identifying disease genes, and finding new drug targets. The discovery of essential genes via wet-lab experimental methods are often time-consuming, laborious, and costly. With the accumulation of gene essentiality data in some model organisms and human cell lines, many computational methods have been proposed to predict essential genes by exploring the correlations between gene essentiality and all sorts of biological information.

One focus in this direction is network based centrality measures. Many studies have demonstrated that highly connected proteins in a protein-protein interaction (PPI) network are more likely to be essential than those of low connections. Although the so-called centrality-lethality rule has been observed in several species, the prediction accuracy is very low for predicting gene essentiality solely based on each of these network topological features. One reason is that the existing PPI networks are not complete and very noisy. The other reason might be the fact that gene essentiality is expected to be affected by multiple biological factors which cannot be fully captured by network topological features. To improve the prediction accuracy, several new centrality measures have been proposed by combining topological properties with other biological information. For example, CoEWC integrated network topological property with gene expression data to capture the common features of essential proteins in both date hubs and party hubs, and showed significant performance improvement compared to methods only based on PPI networks [1]. Zhang et al. proposed an ensemble framework based on gene expression data and PPI networks, which can significantly improve the prediction accuracy of common used centrality measures [2]. Zhang et al. also proposed an integrated method, OGN, by combining network topological properties, the probability of co-expression with the neighboring proteins, and the orthologs in reference organisms [3]. Li et al. proposed GOS [4] by integrating gene expression, orthology, subcellular localization and PPI networks to predict gene essentiality. UDoNC combined the domain features with the topological properties of

PPI networks to predict protein essentiality [5]. Centrality measure based methods predict gene essentiality by a scalar score derived whether from biological network or by integrating multiple data sources, which have limited power for accurately identifying all essential genes. More details about centrality measures for predicting essential genes/proteins can be found in a recent review [6].

The other focus is using machine learning to integrate multiple features derived from different biological data sources to predict gene essentiality. Zhang et al. provided a comprehensive review for gene essentiality predicting methods based on machine learning and network topological features, and pointed out the challenges and potential research directions [7]. As shown in [7], most of the proposed machine learning based predicting methods were evaluated in model organisms. In addition, the traditional machine learning methods were used to predict gene essentiality. In traditional machine learning based essential gene prediction methods, features are often selected and extracted manually, which requires researchers to have prior domain knowledge and keen insights of the relationship between gene essentiality and sorts of biological data in order to obtain the informative features to train the models. For example, Guo et al. used SVM (Support Vector Machines) to predict human gene essentiality based on the  $\lambda$ -interval Z curve derived features from nucleotide sequence data [8]. One limitation of human extracted features would be the coverage. For example, many topological features are derived from PPI network, such as degree centrality, betweenness centrality, closeness centrality, subgraph centrality, and eigen vector centrality, to name a few. The relationship to gene essentiality for each of them has been evaluated by many researchers on several organisms. They have also been used in many machine learning based essential gene prediction methods as shown in [7]. However, their prediction powers either alone or in integration mode (in machine learning methods) are still limited compared with those automatically learned by some deep learning frameworks [9].

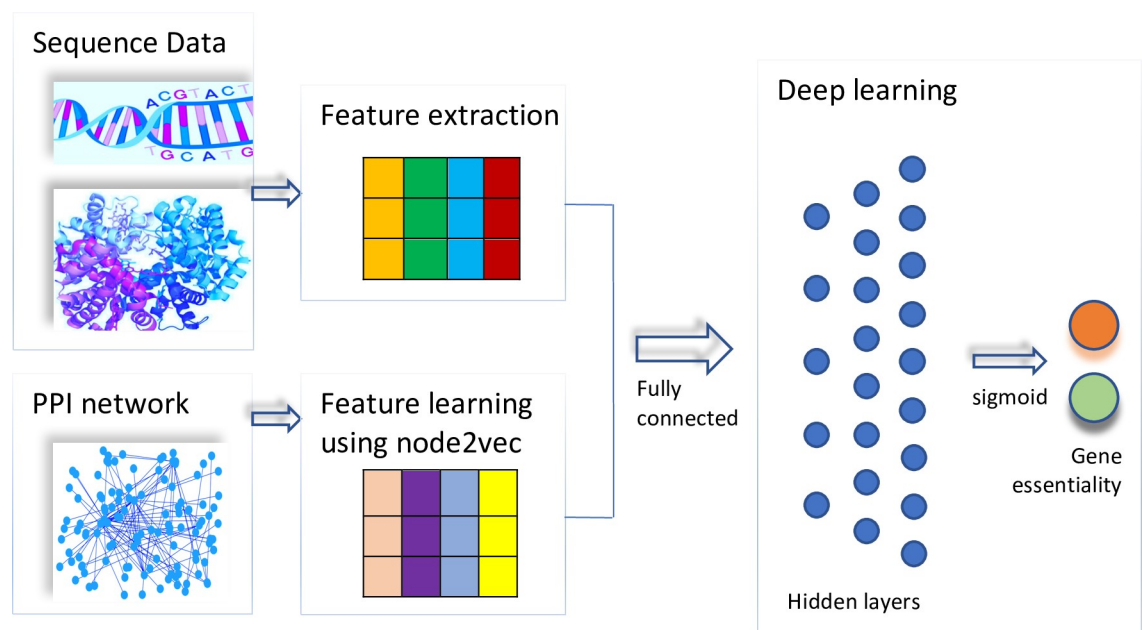
In recent years, deep learning has been applied successfully in many bioinformatics fields, such as medical image segmentation [10], drug-target prediction [11], and essential gene prediction [9, 12, 13]. Convolutional neural network (CNN) has been used to automatically extract features from images [10] or sequence data [11]. Zeng et al. used CNN to extract useful patterns from the time-series gene expression profiles by converting it to an image format based on the cell cycles [9]. Zeng et al. also used bidirectional long short term memory (LSTM) cells to extract features from the same time-series gene expression profiles in their deep learning framework for predicting gene essentiality by integrating gene expression data, subcellular localization data, and PPI networks together, and tested it on *S. cerevisiae* [12]. Hasan et al. used a six hidden-layers neural network to predict gene essentiality in microbes based on sequence data in which features are extracted manually [13]. In addition, deep learning based network embedding methods have been proposed to learn a lower dimensional representation for each node automatically [14]. For example, Zeng et al. used node2vec [14], a deep learning based network embedding method, to extract network features for each protein in a PPI network and showed that such low-dimension representation is more informative than those manually extracted centrality measures [9, 12].

Recently, human essential genes have been identified in several human cancer cell lines using CRISPR-Cas9 and gene-trap technology [15–17]. These identified essential genes provided a clear definition of the requirements for sustaining the basic cell activities of individual human tumor cell types, and can be regarded as targets for cancer treatment [18]. These essential gene datasets together with other available biological data sources enable us to test one important and interesting assumption that human gene essentiality might be accurately predicted using computational methods. Although many previous studies showed that features derived from experimental omics data are useful to predict gene essentiality, such experimental

omics data are often unavailable for under studied organisms. That's why some researchers have been trying to predict essential genes by only utilizing features extracted from sequence data like DNA sequences and protein sequences [8, 13]. In this paper, we propose a deep learning based method to predict human gene essentiality by using features derived from DNA and protein sequence data, which is therefore easily ready to be used for predicting essential genes in other organisms. In addition, in order to improve the performance of the proposed method, we also explore features automatically learned by using a deep learning embedding method from human protein interaction network. We show that each of the two types of features can train a classifier with acceptable prediction performance based on the proposed multiple-layers neural network, and the integration of these features further improves the prediction accuracy.

## Methods

Fig 1 gives the overall architecture of the proposed deep learning framework, DeepHE. It mainly consists of two parts, feature extraction and learning part and classification part. It takes two types of data as input, the sequence data and PPI network. At the feature extraction level, several sequence features for each gene are extracted from its nucleotide sequence and protein sequence data. In addition, an embedding method, node2vec [14], is used to learn the semantic features for each gene from the PPI network. These two types of feature vectors for each gene are concatenated together as the input to the classification module. The classification module is based on multiple-layers neural network which consists of several fully connected hidden layers and an output layer. All its hidden layers utilize the excellent activation function for deep learning, ReLU (Rectified Linear Unit), and use dropout technique to prevent overfitting. After the hidden layers, the fully connected output layer uses sigmoid as its activation function. Considering the skewed distribution nature of human essential gene prediction problem, we explore a cost-sensitive technique to address the imbalanced learning issue when training the classifier by using class weight.



**Fig 1. The flowchart of DeepHE.**

<https://doi.org/10.1371/journal.pcbi.1008229.g001>

## Features derived from sequence data

We extracted features from gene nucleotide sequences and protein sequences. Several features derived from sequence data have been validated their usefulness in predicting gene essentiality in model organisms [13, 19]. In this paper, we used the following sequence derived features: codon frequency, maximum relative synonymous codon usage ( $RSCU_{max}$ ), codon adaptation index (CAI), gene length, GC content, amino acid frequency, and protein sequence length.

Codon frequency of a gene is computed by sliding a window of three nucleotides along its DNA sequence. The raw counts of 64 codons for each gene are calculated and normalized. Unbalanced synonymous codon usage is prevalent in both prokaryotes and eukaryotes. Codon usage bias in a gene may imply its foreign origin, different functional constraints or a different regional mutation. RSCU is a simple measure of non-uniform usage of synonymous codons in a coding sequence, which is defined as the number of times a particular codon is observed, relative to the number of times that the codon would be observed for a uniform synonymous codon usage. Given a synonymous codon  $i$  that has an  $n$ -fold degenerate amino acid, RSCU is computed as (1), where  $X_i$  is the number of occurrence of codon  $i$ , and  $n$  is 1, 2, 3, 4, or 6 according to the genetic code. In this paper, we use the maximal RSCU of each gene as a feature.

$$RSCU_i = \frac{X_i}{\frac{1}{n} \sum_{i=1}^n X_i} \quad (1)$$

Codon adaptation index (CAI) estimates the bias towards certain codon that are more common in highly expressed genes. The CAI of a gene is defined as (2) where  $L$  is the number of codons in the gene excluding methionine, tryptophan, and stop codon.

$$CAI = (\prod_{i=1}^L r_i)^{1/L}, r_i = \frac{RSCU}{RSCU_{max}} \quad (2)$$

In addition to the 68 features derived from gene nucleotide sequences (64 codon frequency and 1 GC content, gene length, CAI, and  $RSCU_{max}$ , respectively), we also use amino acids frequencies and the protein length, that is, 21 features derived from protein sequences. All features are scaled to have mean  $m = 0$  and standard deviation  $std = 1$ .

## Features learned from PPI network

Network embedding methods aim at learning low-dimensional latent representation of nodes in a network, and these representations can be used as features for classification task. Different from some common used topological features, such as node degree centrality (DC), betweenness centrality (BC), and closeness centrality (CC), which usually capture one type of network topological characteristics, the feature representations learned by embedding methods are expected to capture the similarity between nodes in a network.

In this paper, we use a network embedding method, node2vec [14], to automatically learn features for each gene from PPI network. It utilizes a flexible notion of a node's network neighborhood and a biased random walk procedure to learn richer representations. It aims to learn a mapping of nodes to a low-dimensional space of features that maximizes the likelihood of preserving network neighborhoods of nodes. The biased random walk procedure will generate a corpus which consists of many routes each including multiple nodes. These routes just like the sentences including multiple words in natural language. Then these routes will be fed to word2vec framework using a skip-gram technique to learn low-dimensional features for each node. Each node in a PPI network represents a gene/protein. In this way, we get 64 features for each gene from the PPI network.

## Deep learning model based on multilayer perceptron

The classification module in our deep learning framework, DeepHE, is based on the multilayer perceptron structure. It includes one input layer, three hidden layers, and one output layer. All the hidden layers utilize the rectified linear unit (ReLU) activation function. A ReLU is simply defined as  $f(x) = \max(0, x)$ , which turns negative values to zero and grows linearly for positive values. In DeepHE, the output layer uses sigmoid activation function to perform discrete classification. The loss function in DeepHE is binary cross-entropy.

After each hidden layer, a dropout layer is used to make the network less sensitive to noise in the training data and increase its generalization ability. The dropout layer randomly assigns zero weights to a fraction of the neurons in the network. [Table 1](#) gives the parameters used in DeepHE.

The output  $y$  of layer  $i$  depends on the input of layer  $i-1$  as shown in (3), where  $x$  is the input,  $\sigma$  is the activation function,  $b$  is the bias, and  $W$  is the edge weight matrix. During the training phase, the network learns the weights  $W$  and the bias  $b$ .

$$y = \sigma(W^i x^{i-1} + b^{i-1}) \quad (3)$$

In order to tackle the imbalanced classification problem, we used class weight to train a weighted neural network or cost-sensitive neural network. In the weighted neural network, the backpropagation algorithm will be updated to weigh misclassification errors in proportion to the importance of the class. This will allow the model to pay more attention to examples from the minority class than the majority class in datasets with a severely skewed class distribution.

## Results and discussion

### Data collection

DEG database [20] contains 16 human essential gene datasets, among which 13 datasets are from [15–17], and the other three datasets are from [21–23]. We downloaded all the 16 human essential gene datasets for analysis. In total 8,256 human genes are annotated to be essential in at least one of the 16 datasets. [Fig 2](#) shows the distribution of these essential genes across the datasets. According to the assumption that about 10% human genes might be essential genes [16], we select the genes contained at least in 5 datasets as our essential gene dataset, which has 2,024 genes accounting for ~10% of human genes. The DNA sequences and protein sequences for essential genes were downloaded from DEG. We downloaded the genome DNA sequences and protein sequences for all annotated genes from Ensembl [24] (release 97, July 2019). Excluding the 8,256 annotated essential genes in DEG, the other annotated protein coding genes form our nonessential gene dataset, which has 12,697 genes.

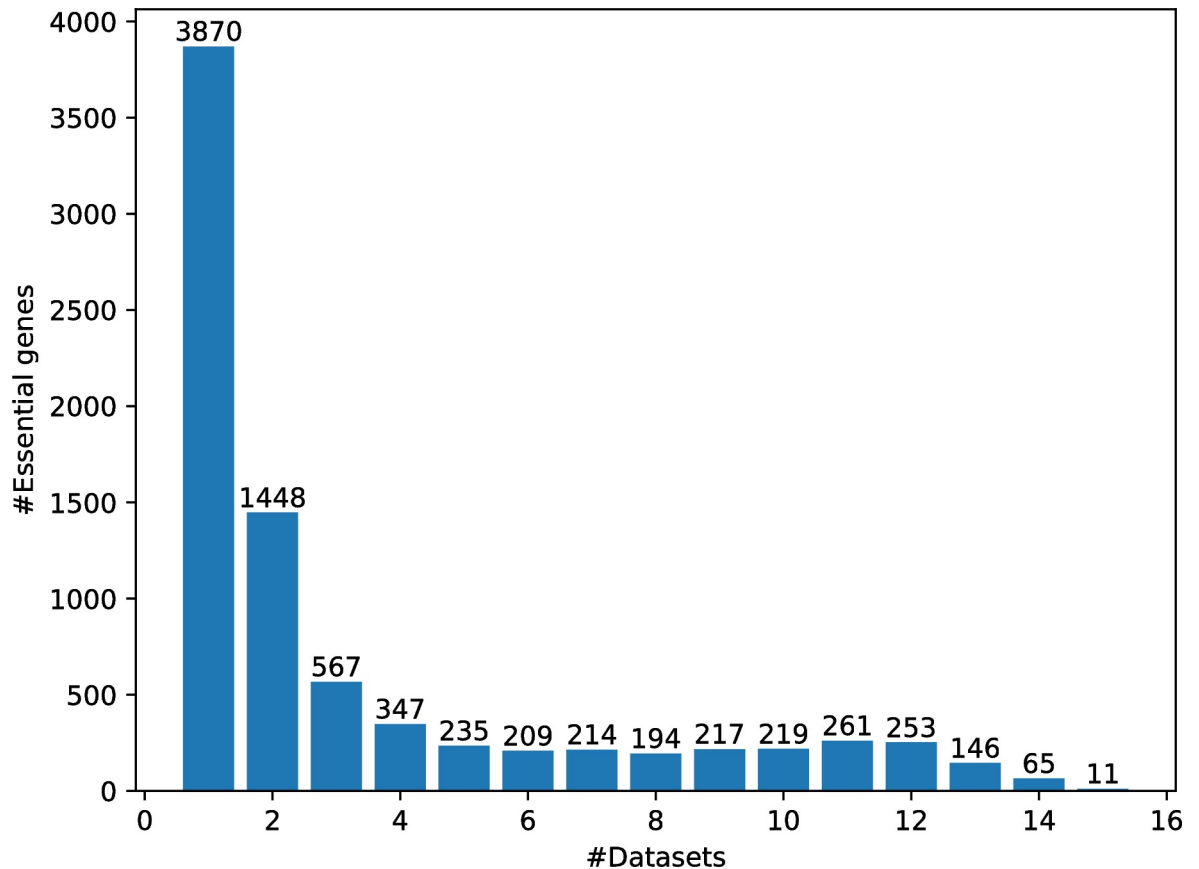
The protein-protein interaction data was downloaded from BioGRID [25] (release 3.5.181, February 2020). Only physical interactions between human genes are used. After filtering out

**Table 1. Parameters of DeepHE.**

	#nodes	Activation function	Dropout probability
Input layer	153	-	-
Hidden layer 1	128	ReLU	0.2
Hidden layer 2	256	ReLU	0.2
Hidden layer 3	512	ReLU	0.2
Output layer	2	sigmoid	-
epochs	100 (early stopping)		
optimizer	Adam (learning_rate = 0.001)		

<https://doi.org/10.1371/journal.pcbi.1008229.t001>





**Fig 2. The distribution of essential genes across the 16 datasets.**

<https://doi.org/10.1371/journal.pcbi.1008229.g002>

self-interactions and several small separated subgraphs, we obtain a protein-protein interaction network with 17,762 nodes and 355,647 edges. This interaction network is used to learn embedding features for each gene. We use genes having both sequence features and network embedding features for training and testing the classification model, that is, 2,009 essential genes and 8,430 nonessential genes are used in the following classification performance evaluations.

The number of nonessential genes is more than 4 folds of that of essential genes, which would suffer the class imbalance problem and result in low predictive accuracy issue for the infrequent class. To address this imbalance issue, class weight is used to train a weighted neural network. In each experiment, the 2009 essential genes and 2009 \* 4 random selected nonessential genes are used to train, validate and test the model. The class weight is set to 4 for the class of essential genes, and 1 for that of nonessential genes. We will also test the effect of different weights to the performance of our model.

### Evaluation metrics

The performance of DeepHE is evaluated using the area under the receiver operating characteristic (ROC) curve (AUC). ROC plot represents the trade-off between sensitivity and specificity for all possible thresholds. We also use the area under the precision-recall curve (AP) to evaluate its performance. Precision-Recall (PR) curves summarize the trade-off between the true positive rate and the positive predictive value for DeepHE using different probability

thresholds. ROC curves are appropriate for balanced classification problems in which each class has almost identical number of instances while PR curves are more appropriate for imbalanced datasets. Since human essential gene prediction is an imbalanced classification problem, the area under the PR curve (AP) should be more indicative than AUC-ROC. In addition to AUC and AP scores, we also give the following performance measures: sensitivity ( $S_n$ ), specificity ( $S_p$ ), positive predictive value (PPV), and accuracy (Ac), which are defined in (4)–(7), where  $TP$ ,  $TN$ ,  $FP$ , and  $FN$  are the number of true positives, true negatives, false positives, and false negatives, respectively.

$$S_n = \frac{TP}{TP + FN} \quad (4)$$

$$S_p = \frac{TN}{FP + TN} \quad (5)$$

$$PPV = \frac{TP}{TP + FP} \quad (6)$$

$$Ac = \frac{TP + TN}{TP + FN + TN + FP} \quad (7)$$

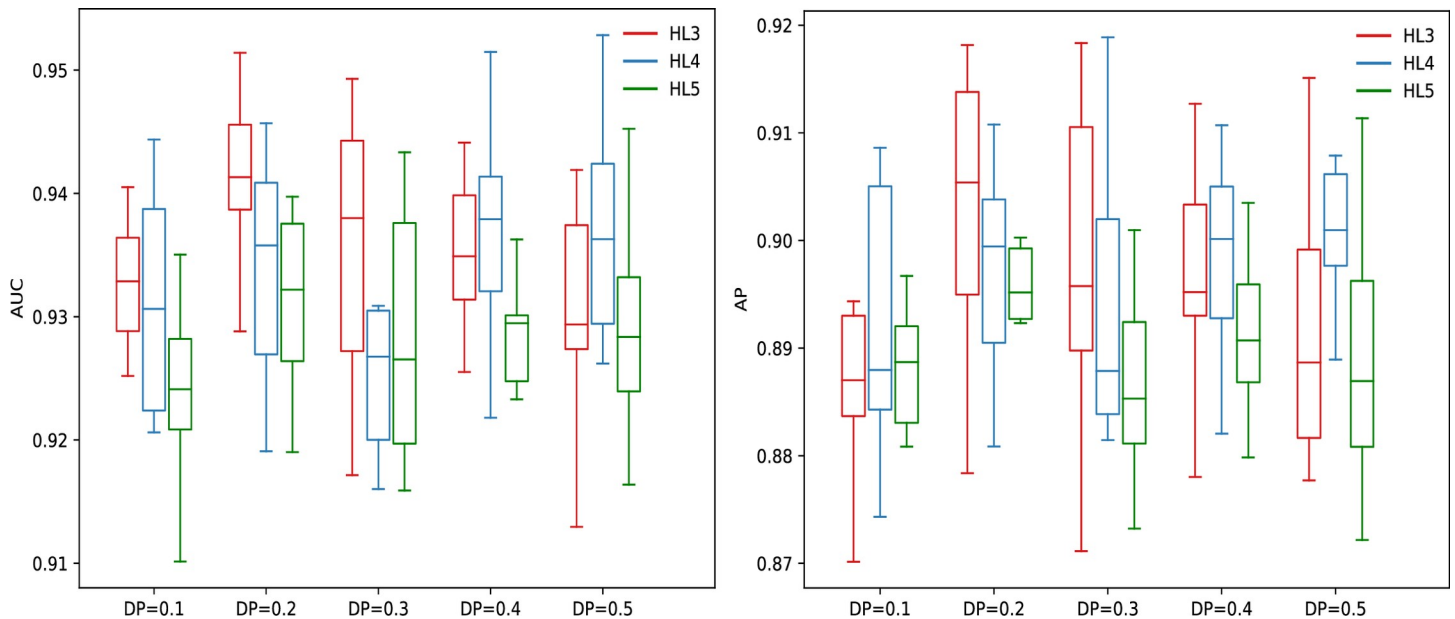
## Performance evaluation

**The effect of number of hidden layers and dropout probability.** There are several hyper-parameters in DeepHE, which would affect its performance. In the following experiments, we choose Adam as the optimizer because of its superior performance. Its initial learning rate is 0.001. The training runs for 100 epochs with early stopping criteria. The batch size is 32. For each run, the 2009 essential genes together with 2009 \* 4 nonessential genes which are randomly selected from the 8430 nonessential genes are used to train, validate, and test the model. We use 80% data for training, 10% data for validation, and the other holding out 10% data for testing. We keep the same ratio between the number of essential genes and that of nonessential genes in training, validating, and testing data. Each experiment is executed 10 times to get the average performance aka 10-fold cross validation with independent 10% data as testing data for each run. The average performance on the holding out independent testing datasets is reported.

Fig 3 and S1 Fig show the performance comparison of DeepHE with different number of hidden layers and different dropout probability (DP). From Fig 3 and S1 Fig we can see that the overall performance of DeepHE is very robust to these two parameters. For example, its best, average, and worst AUC scores are 94.15%, 93.23%, and 92.47% respectively when dropout probability takes values from 0.1 to 0.5 and the number of hidden layers takes values from 3 to 5. It achieves the best overall performance with AUC = 94.15% when using HL3 with DP = 0.2. Its AP scores are also very stable with the best, average, and worst values of 90.64%, 89.4%, and 88.69% respectively. Same with AUC, it achieves the best AP score of 90.64% when using HL3 with DP = 0.2. In addition to the best AUC and AP scores, it also achieves the best scores for specificity (94.5%), PPV (77.74%), and accuracy (90.88%) when using HL3 with DP = 0.2. The best sensitivity score is 87.16% when using HL5 with DP = 0.5. From S1 Fig we can also see that with the increase of dropout probability, its sensitivity score increases but its PPV score decreases in most cases.

In a very skewed classification problem, the accuracy and AUC measures can get large values even when almost all the instances in the minority class are classified into the majority





**Fig 3. Performance comparison of DeepHE with different number of hidden layers and different dropout probability (DP) for two metrics: AUC and AP.** HL3 = [128, 256, 512], HL4 = [128, 256, 512, 1024], HL5 = [128, 256, 512, 1024, 1024]. DP: dropout probability.

<https://doi.org/10.1371/journal.pcbi.1008229.g003>

class. That's not what we expected. In most cases of imbalanced classification problems, we are far more concerned with the classifier's performance on the minority class. Since the essential gene prediction problem is often a very skewed classification problem in which the number of essential genes is much less than that of nonessential genes. Our concerns would be how many essential genes can be predicted and how many genes are truly essential among those predicted as essential genes, that is, sensitivity and PPV as well as the comprehensive measure AP are more important. Based on this point, we think that DeepHE with 3 hidden layers and DP = 0.2 is the best one which will be used in the following experiments. Although the other parameters in DeepHE would also impact its performance, we only test the impact of dropout probability and the number of hidden layers while keeping other parameters to fixed values since our aim is not to find out the best parameters here. It's possible that DeepHE performs even better after fine tuning the other parameters.

Fig 4 gives the ROC curves of DeepHE in 10 repetitions when using HL3 and DP = 0.2. ROC curves summarize the trade-off between the true positive rate and false positive rate of DeepHE using different probability thresholds. From Fig 4 we can see that DeepHE reached its best performance at iterations 2, 4, 8, and 10 with AUC = 0.95. In addition, the performance of DeepHE is quite stable since the difference is only about 0.02 between its best and worst AUC scores. Guo et al. also used machine learning (SVM) to predict human essential genes based on sequence data [8]. Their reported best performance is AUC = 0.88. Compared with [8], DeepHE outperforms their method.

Fig 4 also shows PR curves for 10 iterations of DeepHE with HL3 and DP = 0.2. Similar with the AUC scores, its AP scores are also very stable since there's only a very small difference between its best and worst AP scores (about 0.04). It achieves the best performance in iteration 6 with AP = 93%. The worst AP score is still above 88% which indicates that DeepHE is very effective for predicting human essential genes.

**The effect of class weight.** In order to cope with the imbalanced data distributions between two classes, DeepHE uses class weight to give larger penalty when misclassifying an

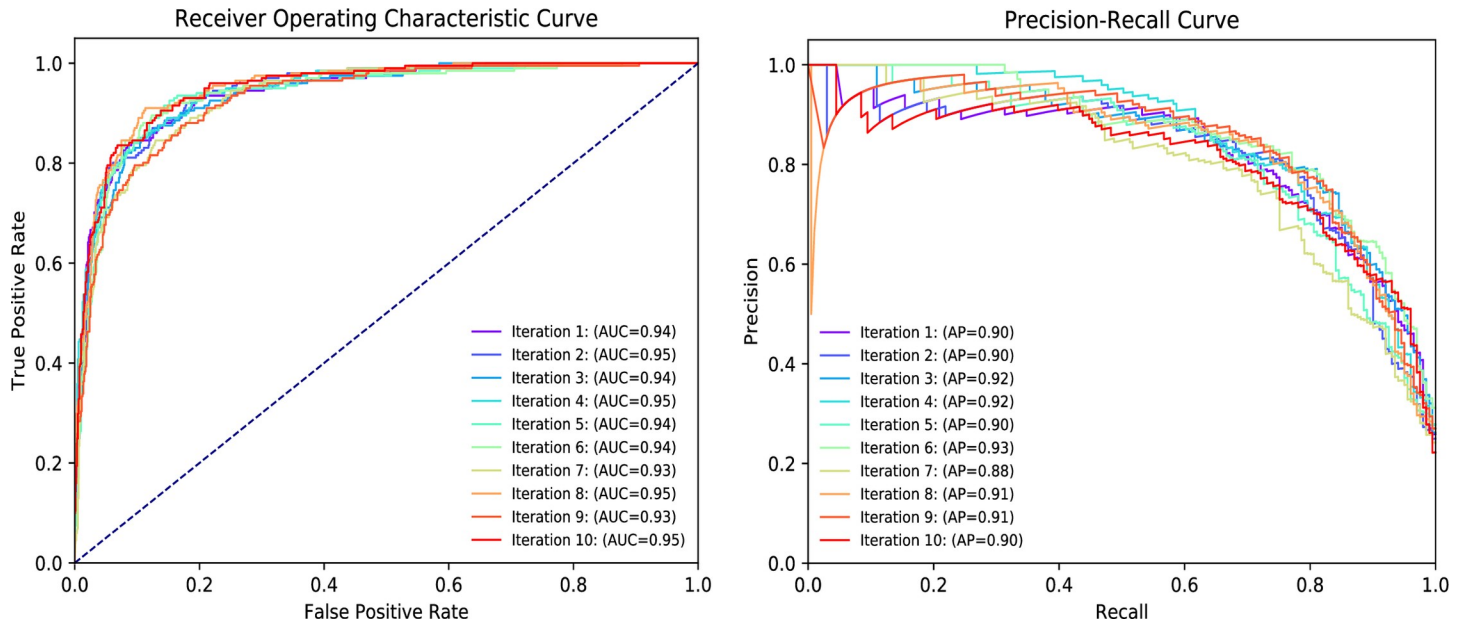


Fig 4. The ROC and PR curves of DeepHE with HL3 and DP = 0.2.

<https://doi.org/10.1371/journal.pcbi.1008229.g004>

instance in the minority class, that is, the class of essential genes. In the following, we will test if different class weight values would affect the performance of DeepHE. Note that in each experiment, the ratio between the number of essential genes and that of nonessential genes is 1:4. The class weight for nonessential genes is always 1. We will vary the class weight for essential genes from 1 to 10 to see its effect on the performance. DeepHE with 3 hidden layers and DP = 0.2 is used for the following experiments.

Fig 5 and S2 Fig give the performance comparison of DeepHE with different class weights for the class of essential genes. From Fig 5 and S2 Fig we can see that DeepHE achieves best

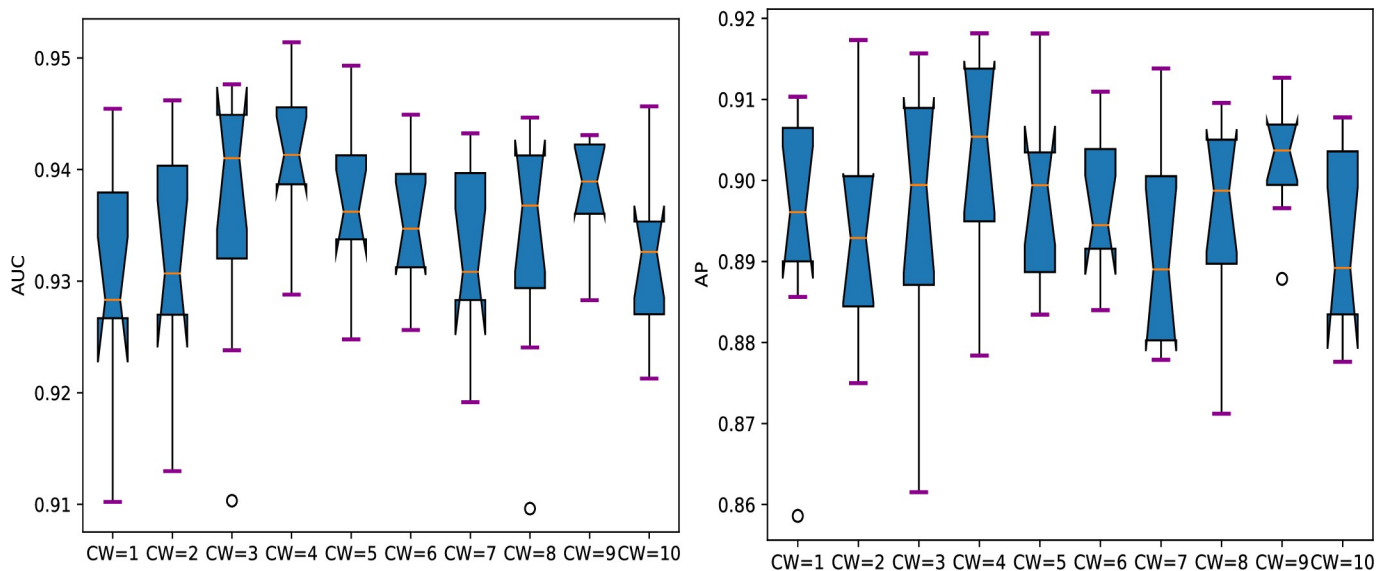


Fig 5. Performance comparison of DeepHE with different class weights for two metrics: AUC and AP. CW: class weight.

<https://doi.org/10.1371/journal.pcbi.1008229.g005>

AUC (94.15%), PPV (77.74%), Accuracy (90.88%), and AP score (90.64%) when class weight equals 4 for essential genes. It gets the best sensitivity score (79.85%) when class weight = 9. In general, the sensitivity score increases with the increase of the class weight, but PPV score decreases with the increase of class weight. This accords with our intuition. With larger class weight for essential genes, misclassifying an essential gene will get larger penalty than misclassifying a nonessential gene. In this situation, more essential genes will be put into the right class, at the same time, more nonessential genes would also be put into the class of essential genes, which will result in higher sensitivity score and lower PPV score. When class weight = 4, it mimics the situation that the number of essential genes equal to the number of nonessential genes, thus it achieves a balanced point for sensitivity and PPV score. One can set the class weight according to his preference to whether higher sensitivity or higher PPV or just the balance between them.

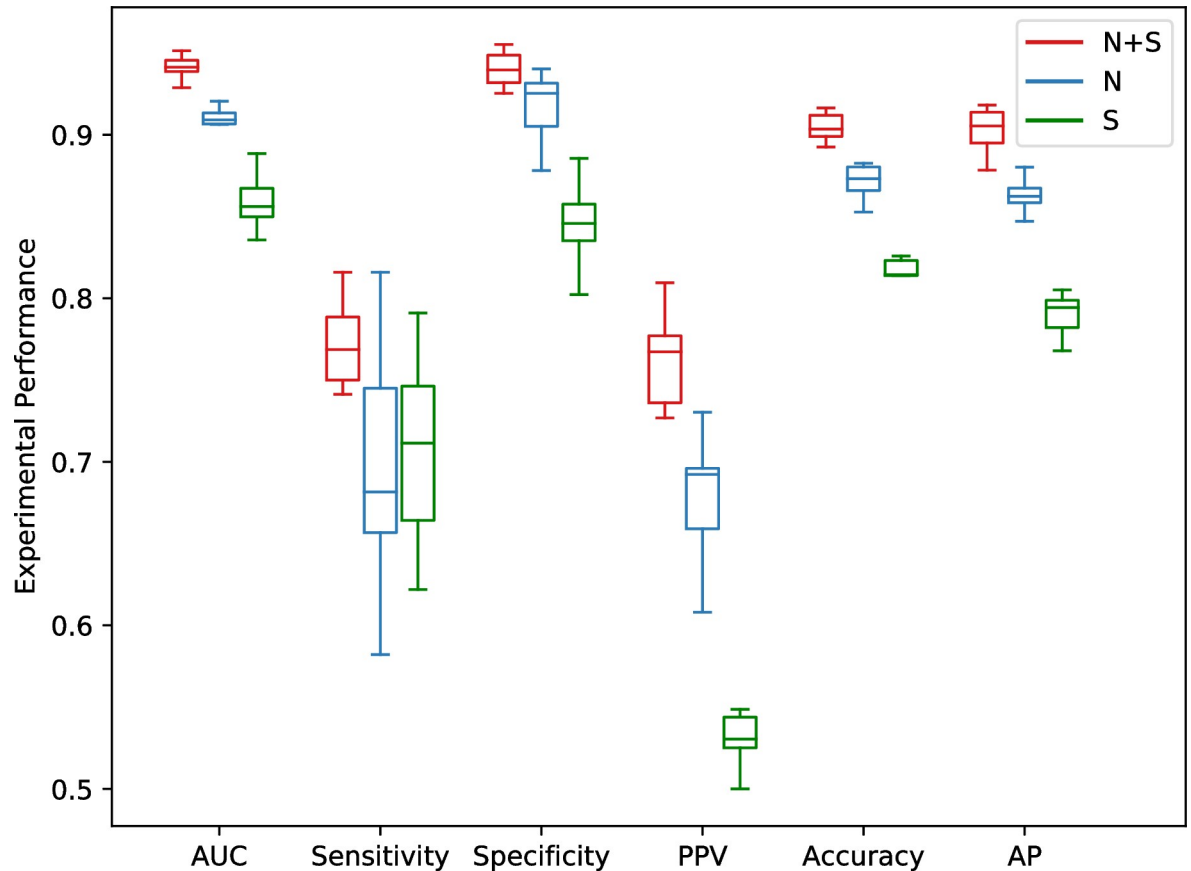
[Fig 5](#) and [S2 Fig](#) also tell us that the AUC, specificity, AP, and Accuracy of DeepHE are very robust to the class weight. For example, the best, average, and worst AUC scores are 94.15%, 93.48%, and 93.07% respectively; the best, average, and worst specificity scores are 94.75%, 93.03%, and 90.95% respectively; the best, average, and worst AP scores are 90.64%, 89.69%, and 89.19% respectively; the best, average, and worst Accuracy scores are 90.88%, 89.69%, and 88.39% respectively. When varying the class weight, sensitivity score and PPV change in opposite directions which makes the overall performance of DeepHE only slightly affected by the change of class weight.

**The contribution of different features.** DeepHE utilizes two types of features, sequence features (S) and network embedding features (N). In the following we will test how each type of features affect the performance of DeepHE. In the following experiments, DeepHE works with same configurations (3 hidden layers, DP = 0.2, class weight = 4. Other configurations are same as before) except the input features.

[Fig 6](#) gives the performance of DeepHE using different type of features. It tells us that DeepHE with the integration of sequence features and network embedding features works best which confirms the contribution and complement of the two types of features. DeepHE with only sequence features works worst which has very low PPV score (53.28%). DeepHE with network embedding features works in between, whose AP score achieves acceptable level (86.53%). DeepHE achieves the best performance for all the six measures by integrating these two types of features.

**Comparison with centrality measures.** In order to demonstrate the effectiveness of DeepHE, we compare it with several popular centrality measures which are widely used either alone or as features in machine learning based methods to predict essential genes/proteins. Four commonly used centrality measures (DC, BC, EC, CC) are used for the comparison in the following steps. First, each centrality measure is used to compute the values of proteins in the PPI network (the same PPI network is used as for learning embedding features). Second, the proteins are ranked by descending order. Third, the top ranked 2009 genes are selected as candidate essential genes. Based on this partition, we can calculate accuracy, PPV, and sensitivity according to the true labels of genes. [Fig 7](#) shows the performance comparison of these methods. From [Fig 7](#) we can see that DeepHE outperforms the other centrality measures regards to all the three metrics. For example, the accuracy of DeepHE increases about 8.5%, 14%, 5.4%, and 6.6% compared with that of DC, BC, EC, and CC respectively; the PPV of DeepHE increases about 33.7%, 64%, 20.2%, and 15.5% compared with that of DC, BC, EC, and CC respectively. The sensitivity of DeepHE has the similar percent improvement with that of PPV. This demonstrates that DeepHE is superior that the commonly used centrality measures, which accords with our intuition.

In order to test whether the features automatically learned by network embedding method are more informative than that of manually designed centrality measures for essential gene

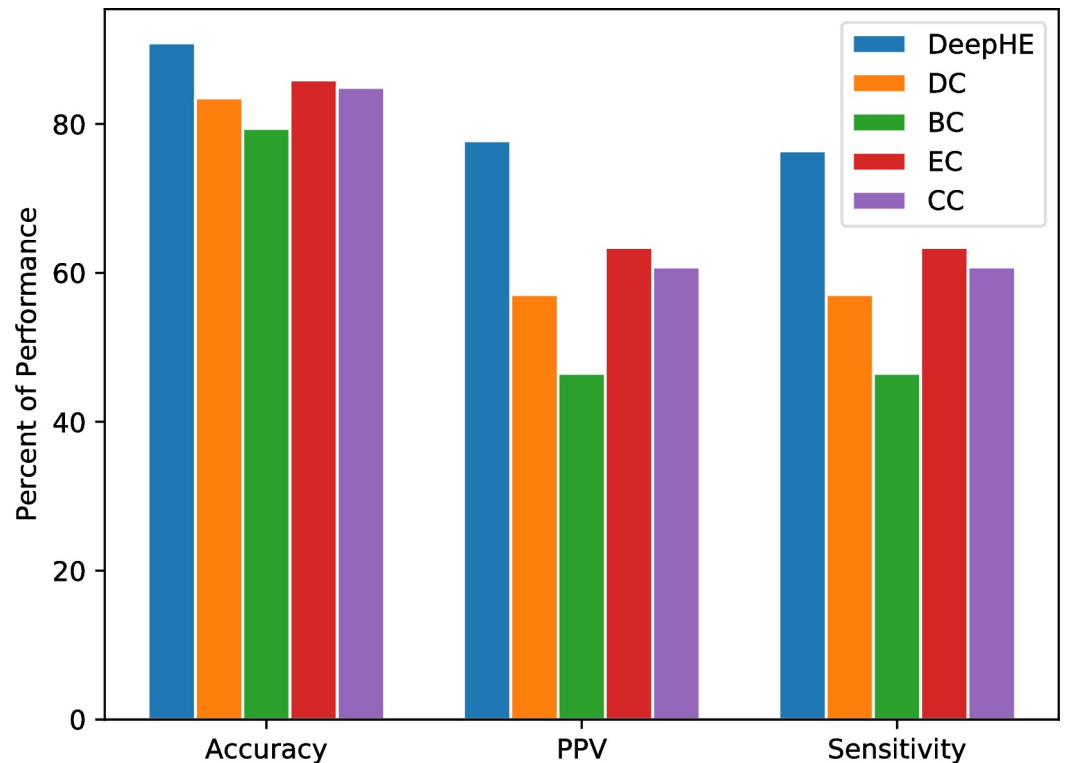


**Fig 6. Performance comparison of DeepHE with different features.** N+S: network embedding features plus sequence features; N: network embedding features; S: sequence features.

<https://doi.org/10.1371/journal.pcbi.1008229.g006>

prediction, we further run several experiments for DeepHE with different features, in which sequence features are combined with that calculated by one of the four widely used centrality measures. The performance results are shown in Fig 8. It tells us that DeepHE(N+S) outperforms all the other variants by substituting network embedding features with that from one of the centrality measures regards to the metrics of AUC, AP, accuracy, PPV, and specificity. For instance, compared with that of DeepHE using S+DC, S+BC, S+EC, and S+CC as features respectively, the average AUC of DeepHE(N+S) increases about 3.3%, 7%, 1.8%, and 3.2%; the average AP of DeepHE(N+S) increases about 4.4%, 8.9%, 2.5%, and 4.1%; the average accuracy of DeepHE(N+S) increases about 5.8%, 6.3%, 3.8%, and 5.2%. From Fig 8 we can also see that DeepHE with S+BC performs worst for almost all the metrics, which tells us that betweenness centrality captures less useful information from the PPI network for the task of essential gene prediction. This can also be confirmed in Fig 7. BC itself performs worst (Fig 7) so that it cannot provide more complementary information to enhance sequence features.

**Comparison with traditional machine learning models.** Several machine learning methods have been used to predict essential genes [7]. In order to demonstrate the superior of our proposed prediction method DeepHE, we also compare it with several widely used traditional machine learning models, such as Support Vector Machines (SVM), Naïve Bayes (NB), Random Forest (RF), and Adaboost. All the compared machine learning algorithms are implemented by scikit-learn python library with default parameters, unless otherwise specified. For each model, we either set class weight parameter to 4 or set sample weight parameter to 4 for



**Fig 7. Performance comparison of DeepHE, DC, BC, EC, and CC.**

<https://doi.org/10.1371/journal.pcbi.1008229.g007>

each essential gene and 1 for each nonessential gene, therefore the two types of weights are essentially same. The sample weight is only used when class weight is not available. All models are tested 10 times and the average performance for each measure is used for comparison. More performance information about these tests can be found from the box plots in Fig 9 and S3 Fig.

From Fig 9 and S3 Fig we can see that DeepHE (N+S) significantly outperforms the other machine learning models regarding to three comprehensive measures, AUC, AP, and Accuracy. For instance, the AUC score of DeepHE (N+S) is 8.79% higher than that of SVM (N+S), 43.22% higher than that of NB (N+S), 25.38% higher than that of RF (N+S), and 15.39% higher than that of Adaboost (N+S). The AP score of DeepHE (N+S) increases by 51.24%, 220.77%, 65.95%, and 91.63% compared with that of SVM (N+S), NB (N+S), RF (N+S), and Adaboost (N+S) respectively. By integrating sequence features and network embedding features, the overall performance of four models (DeepHE, SVM, RF, Adaboost) gets improved. NB works slightly better with only network embedding features. Considering the fact that essential gene prediction is an imbalanced problem, AP is more important than other measures. From Fig 9 we can see that the four compared traditional machine learning models have very low AP scores (from 27.91% to 59.93%), which tells us that they are not a good choice for such task, and further confirms the superior of our proposed deep learning model, DeepHE.

## Conclusion

We propose a new essential gene prediction framework based on deep learning, DeepHE. It aims to explore whether deep learning can achieve notable improvements for predicting gene essentiality, an imbalanced classification problem. DeepHE integrates two types of features,

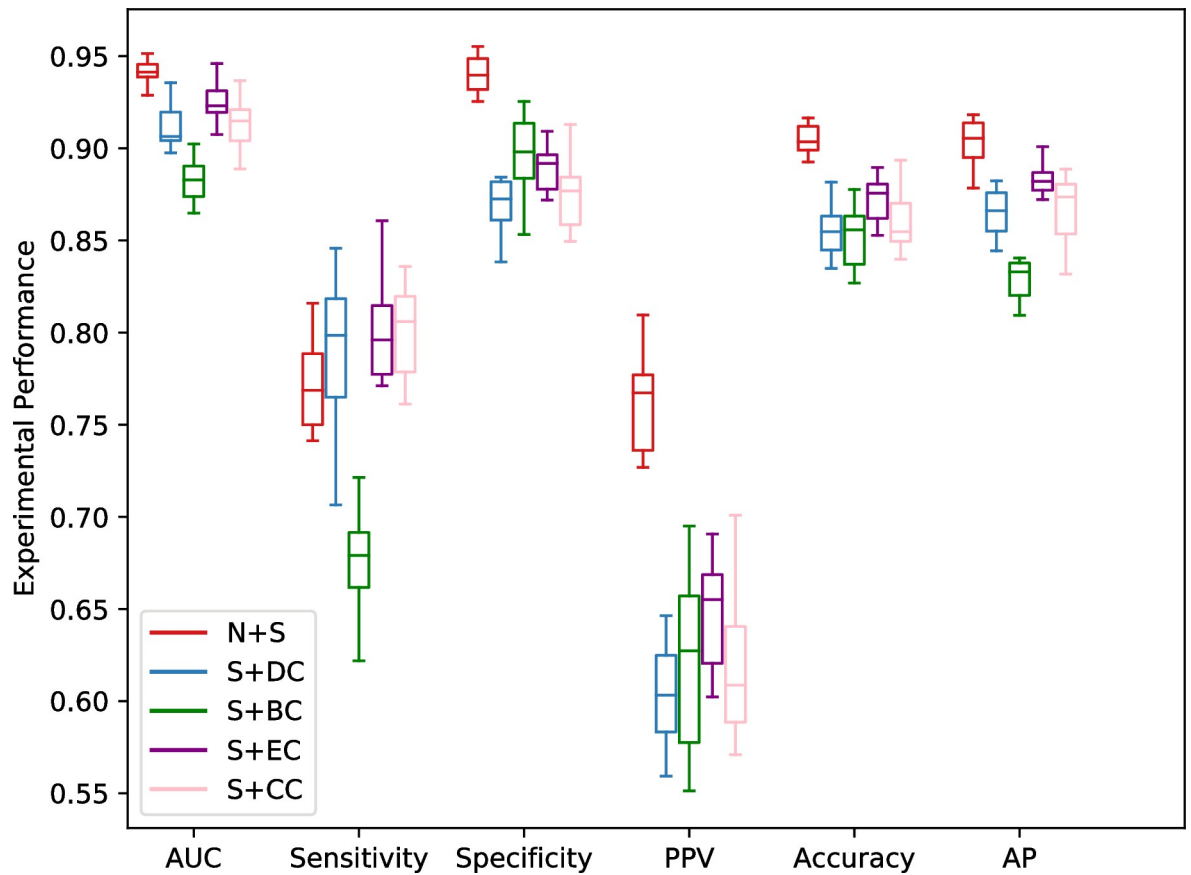


Fig 8. Performance comparison of DeepHE with different feature combinations.

<https://doi.org/10.1371/journal.pcbi.1008229.g008>

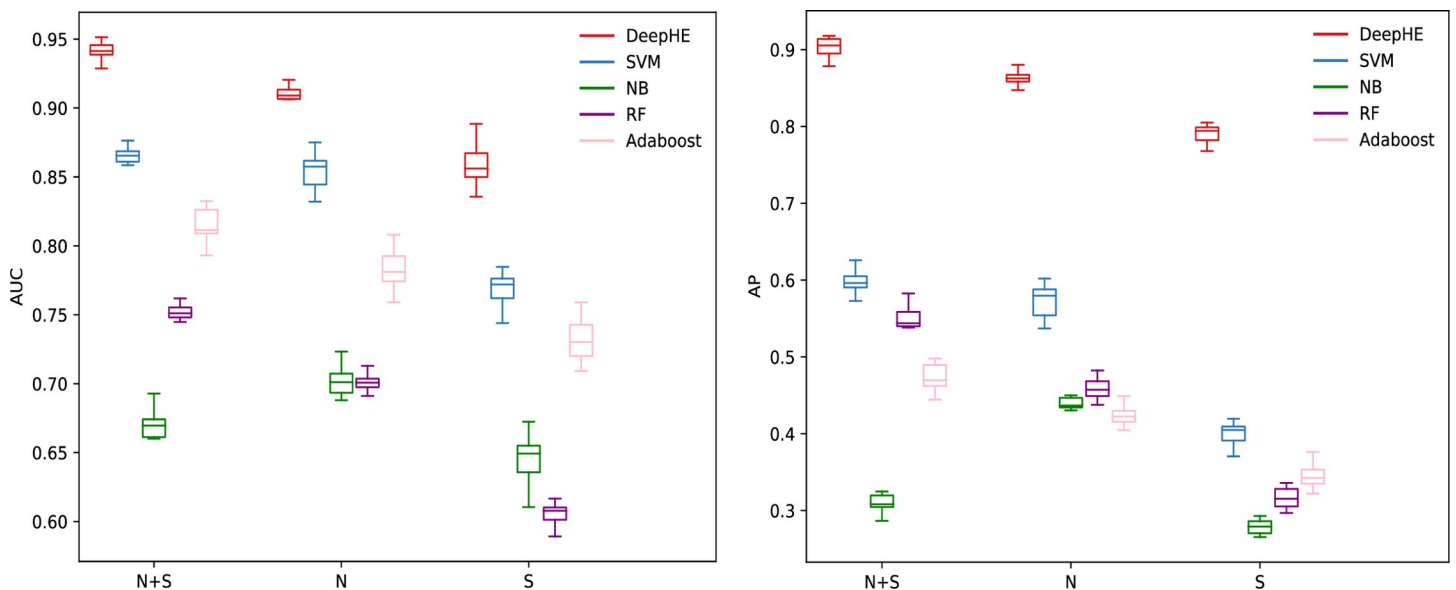


Fig 9. Performance comparison of DeepHE, SVM, NB, RF, and Adaboost with different features. N: network embedding features; S: sequence features.

<https://doi.org/10.1371/journal.pcbi.1008229.g009>



sequence features extracted from DNA sequence and protein sequence and features learned from PPI network, as its input. Then a multilayer perceptron was used to train a cost-sensitive classifier by setting class weight. Although several machine learning based essential gene prediction methods have been proposed, most of them base on the features extracted according to human domain knowledge. In this paper, we use a deep learning based network embedding algorithm, node2vec, to automatically learn network features for each gene from the PPI network. The learned embedding features greatly improved the performance of DeepHE compared with it either only using sequence features or using sequence features plus one feature computed by one of widely used centrality measures. The performance of DeepHE is evaluated on human datasets, which achieves very good performance for three comprehensive measures AUC (94.15%), AP (90.64%), and Accuracy (90.88%). We also compare it with four widely used machine learning models, SVM, Naïve Bayes, Random Forest, and Adaboost, as well as four popular centrality measures (DC, BC, EC, and CC). DeepHE significantly outperforms all the four machine learning models and the four centrality measures, which further demonstrates that DeepHE is an effective deep learning framework for human essential gene prediction.

In the future, we will explore and integrate other biological data to further improve the performance of DeepHE. Especially we are interested in how to use deep learning to automatically learn features from biological data rather than manually extracting features heavily based on domain knowledge. In addition, we are also interested in exploring more useful techniques to cope with the imbalanced classification problem as well as sparsely labeled classification problem [26,27], and utilizing membrane computing techniques [28] to enhance the learning procedure. Exploring deep learning to predict human essential genes across human cancer cell lines would be also interesting.

## Supporting information

**S1 Fig. Performance comparison of DeepHE with different hidden layers and dropout probability for four metrics: accuracy, PPV, sensitivity, and specificity.**

(TIF)

**S2 Fig. Performance comparison of DeepHE with different class weights for four metrics: accuracy, PPV, sensitivity, and specificity.**

(TIF)

**S3 Fig. Performance comparison of DeepHE, SVM, NB, RF, and Adaboost with different features for four metrics: accuracy, PPV, sensitivity, and specificity.**

(TIF)

## Author Contributions

**Conceptualization:** Xue Zhang.

**Data curation:** Xue Zhang.

**Funding acquisition:** Xue Zhang, Wangxin Xiao.

**Investigation:** Xue Zhang, Wangxin Xiao.

**Methodology:** Xue Zhang.

**Software:** Xue Zhang, Weijia Xiao.

**Supervision:** Xue Zhang, Wangxin Xiao.

**Validation:** Xue Zhang, Weijia Xiao.

**Visualization:** Xue Zhang, Weijia Xiao.

**Writing – original draft:** Xue Zhang, Wangxin Xiao.

**Writing – review & editing:** Xue Zhang, Wangxin Xiao.

## References

1. Zhang X, Xu J, Xiao W. A new method for the discovery of essential proteins. *PLoS ONE*. 2013; 8 (3): e58763. <https://doi.org/10.1371/journal.pone.0058763> PMID: 23555595
2. Zhang X, Xiao W, Acencio ML, Lemke N, Wang X. An ensemble framework for identifying essential proteins. *BMC Bioinformatics*. 2016; 17:322. <https://doi.org/10.1186/s12859-016-1166-7> PMID: 27557880
3. Zhang X, Xiao W, Hu X. Predicting essential proteins by integrating orthology, gene expressions, and PPI networks. *PLoS ONE*. 2018; 13(4): e0195410. <https://doi.org/10.1371/journal.pone.0195410> PMID: 29634727
4. Li G, Li M, Wang J, Wu J, Wu F, Pan Y. Predicting essential proteins based on subcellular localization, orthology and PPI networks. *BMC Bioinformatics*. 2016; 17 Suppl 8:279.
5. Peng W, Wang J, Cheng Y, Lu Y, Wu F, Pan Y. UDoNC: an algorithm for identifying essential proteins based on protein domains and protein-protein interaction networks. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. 2015; 12(2): 276–288. <https://doi.org/10.1109/TCBB.2014.2338317> PMID: 26357216
6. Li X, Li W, Zeng M, Zheng R, Li M. Network-based methods for predicting essential genes or proteins: a survey. *Briefings in Bioinformatics*. 2019; bbz017.
7. Zhang X, Acencio ML, Lemke N. Predicting essential genes and proteins based on machine learning and network topological features: a comprehensive review. *Front. Physiol*. 2016; 7:75. <https://doi.org/10.3389/fphys.2016.00075> PMID: 27014079
8. Guo F, Dong C, Hua H, Liu S, Luo H, Zhang H, et al. Accurate prediction of human essential genes using only nucleotide composition and association information. *Bioinformatics*. 2017; 33(12): 1758–1764. <https://doi.org/10.1093/bioinformatics/btx055> PMID: 28158612
9. Zeng M, Li M, Wu F, Li Y, Pan Y. DeepEP: a deep learning framework for identifying essential proteins. *BMC Bioinformatics*. 2019; 20 Suppl 16: 506. <https://doi.org/10.1186/s12859-019-3076-y> PMID: 31787076
10. Ronneberger O, Fischer P, Brox T. U-Net: convolutional networks for biomedical image segmentation. *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, Springer, LNCS. 2015; 9351: 234–241.
11. Ozturk H, Ozgur A, Ozkirimli E. DeepDTA: deep drug-target binding affinity prediction. *Bioinformatics*. 2018; 34, i821–i829.
12. Zeng M, Li M, Fei Z, Wu F, Li Y, Pan Y, et al. A deep learning framework for identifying essential proteins by integrating multiple types of biological information. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. 2019 Feb 5.
13. Hasan MA, Lonardi S. DEEPLYESSENTIAL: A deep neural network for predicting essential genes in microbes. *BioRxiv*. 2019.
14. Grover A, Leskovec J. node2vec: scalable feature learning from networks. *KDD'16: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. August 2016; 855–864.
15. Blomen VA, Májek P, Jae LT, Bigenzahn JW, Nieuwenhuis J, Staring J, et al. Gene essentiality and synthetic lethality in haploid human cells. *Science*. 2015; 350(6264):1092–6. <https://doi.org/10.1126/science.aac7557> PMID: 26472760
16. Wang T, Birsoy K, Hughes NW, Krupczak KM, Post Y, Wei JJ, et al. Identification and characterization of essential genes in the human genome. *Science*. 2015; 350(6264): 1096–101. <https://doi.org/10.1126/science.aac7041> PMID: 26472758
17. Hart T, Chandrashekar M, Aregger M, Steinhart Z, Brown KR, MacLeod G, et al. High-resolution CRISPR screens reveal fitness genes and genotype-specific cancer liabilities. *Cell*. 2015; 163(6): 1515–26. <https://doi.org/10.1016/j.cell.2015.11.015> PMID: 26627737
18. Fraser A. Essential human genes. *Cell Systems*. 2015; 1(6): 381–382. <https://doi.org/10.1016/j.cels.2015.12.007> PMID: 27136352

19. Liu X, Wang BJ, Xu L, Tang HL, Xu GQ. Selection of key sequence-based features for prediction of essential genes in 31 diverse bacterial species. *PLoS ONE*. 2017; 12(3): e0174638. <https://doi.org/10.1371/journal.pone.0174638> PMID: 28358836
20. Luo H, Lin Y, Gao F, Zhang CT, Zhang R. DEG 10, an update of the database of essential genes that includes both protein-coding genes and noncoding genomic elements. *Nucleic Acids Res*. Jan 2014; 42 (Database issue): D574–80. <https://doi.org/10.1093/nar/gkt1131> PMID: 24243843
21. Liao BY, Zhang J. Null mutations in human and mouse orthologs frequently result in different phenotypes. *Proc Natl Acad Sci*. 2008; 105(19): 6987–6992. <https://doi.org/10.1073/pnas.0800387105> PMID: 18458337
22. Georgi B, Voight BF, Bućan M. From mouse to human: evolutionary genomics analysis of human orthologs of essential genes. *PLoS genetics*. 2013; 9(5): e1003484. <https://doi.org/10.1371/journal.pgen.1003484> PMID: 23675308
23. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature*. 2016; 536: 285–291. <https://doi.org/10.1038/nature19057> PMID: 27535533
24. Ruffier M, Kähäri A, Komorowska M, Keenan S, Laird M, et al. Ensembl core software resources: storage and programmatic access for DNA sequence and genome annotation. *Database*, 2017, <https://doi.org/10.1093/database/bax020> PMID: 28365736
25. Stark C, Breitkreutz BJ, Reguly T, Boucher L, Breitkreutz A, Tyers M. Biogrid: a general repository for interaction datasets. *Nucleic Acids Res*. Jan 1, 2006; 34: D535–9. <https://doi.org/10.1093/nar/gkj109> PMID: 16381927
26. Zhang X, Xiao W. Clustering based two-stage text classification requiring minimal training data. *Computer Science and Information Systems*. 2012; 9(4):1627–1643.
27. Zhang X, Xiao W. Active semi-supervised framework with data editing. *Computer Science and Information Systems*. 2012; 9(4): 1513–1532.
28. Zhang G, Pérez-Jiménez MJ, Gheorghe M. Real-life applications with membrane computing. Springer 2017, ISBN 978-3-319-55989-6.