

KG–CNNDTI: a knowledge graph–enhanced prediction model for drug–target interactions and application in virtual screening of natural products against Alzheimer’ s disease

Chengyuan Yue, Baiyu Chen, Long Chen, Le Xiong, Changda Gong, Ze Wang, Guixia Liu, Weihua Li, Rui Wang, Yun Tang

Citation: Chengyuan Yue, Baiyu Chen, Long Chen, Le Xiong, Changda Gong, Ze Wang, Guixia Liu, Weihua Li, Rui Wang, Yun Tang, KG–CNNDTI: a knowledge graph–enhanced prediction model for drug–target interactions and application in virtual screening of natural products against Alzheimer’ s disease, *Chinese Journal of Natural Medicines*, 2025, 23(11), 1283–1292. doi: [10.1016/S1875-5364\(25\)60980-0](https://doi.org/10.1016/S1875-5364(25)60980-0).

View online: [https://doi.org/10.1016/S1875-5364\(25\)60980-0](https://doi.org/10.1016/S1875-5364(25)60980-0)

Related articles that may interest you

[Approved drugs and natural products at clinical stages for treating Alzheimer’ s disease](#)

Chinese Journal of Natural Medicines. 2024, 22(8), 699–710 [https://doi.org/10.1016/S1875-5364\(24\)60606-0](https://doi.org/10.1016/S1875-5364(24)60606-0)

[Progress in approved drugs from natural product resources](#)

Chinese Journal of Natural Medicines. 2024, 22(3), 195–211 [https://doi.org/10.1016/S1875-5364\(24\)60582-0](https://doi.org/10.1016/S1875-5364(24)60582-0)

[Dual–function natural products: Farnesoid X receptor agonist/inflammation inhibitor for metabolic dysfunction–associated steatotic liver disease therapy](#)

Chinese Journal of Natural Medicines. 2024, 22(11), 965–976 [https://doi.org/10.1016/S1875-5364\(24\)60706-5](https://doi.org/10.1016/S1875-5364(24)60706-5)

[Modulation of type I interferon signaling by natural products in the treatment of immune–related diseases](#)

Chinese Journal of Natural Medicines. 2023, 21(1), 3–18 [https://doi.org/10.1016/S1875-5364\(23\)60381-4](https://doi.org/10.1016/S1875-5364(23)60381-4)

[Identification of multi–target anti–cancer agents from TCM formula by *in silico* prediction and *in vitro* validation](#)

Chinese Journal of Natural Medicines. 2022, 20(5), 332–351 [https://doi.org/10.1016/S1875-5364\(22\)60180-8](https://doi.org/10.1016/S1875-5364(22)60180-8)

[Transdermal delivery of natural products against atopic dermatitis](#)

Chinese Journal of Natural Medicines. 2024, 22(12), 1076–1088 [https://doi.org/10.1016/S1875-5364\(24\)60681-3](https://doi.org/10.1016/S1875-5364(24)60681-3)



Wechat



Contents lists available at ScienceDirect

Chinese Journal of Natural Medicines

journal homepage: www.cjnmcpu.com/

Original article

KG-CNNNTI: a knowledge graph-enhanced prediction model for drug-target interactions and application in virtual screening of natural products against Alzheimer's disease



Chengyuan Yue^A, Baiyu Chen^A, Long Chen, Le Xiong, Changda Gong, Ze Wang, Guixia Liu, Weihua Li, Rui Wang*, Yun Tang*

Shanghai Frontiers Science Center of Optogenetic Techniques for Cell Metabolism, Shanghai Key Laboratory of New Drug Design, School of Pharmacy, East China University of Science and Technology, Shanghai 200237, China

ARTICLE INFO

Article history:

Received 24 January 2025

Revised 25 February 2025

Accepted 21 May 2025

Available online 20 November 2025

Keywords:

Drug-target interactions prediction

Knowledge graph

Drug screening

Alzheimer's disease

Natural products

ABSTRACT

Accurate prediction of drug-target interactions (DTIs) plays a pivotal role in drug discovery, facilitating optimization of lead compounds, drug repurposing and elucidation of drug side effects. However, traditional DTI prediction methods are often limited by incomplete biological data and insufficient representation of protein features. In this study, we proposed KG-CNNNTI, a novel knowledge graph-enhanced framework for DTI prediction, which integrates heterogeneous biological information to improve model generalizability and predictive performance. The proposed model utilized protein embeddings derived from a biomedical knowledge graph via the Node2Vec algorithm, which were further enriched with contextualized sequence representations obtained from ProteinBERT. For compound representation, multiple molecular fingerprint schemes alongside the Uni-Mol pre-trained model were evaluated. The fused representations served as inputs to both classical machine learning models and a convolutional neural network-based predictor. Experimental evaluations across benchmark datasets demonstrated that KG-CNNNTI achieved superior performance compared to state-of-the-art methods, particularly in terms of Precision, Recall, F1-Score and area under the precision-recall curve (AUPR). Ablation analysis highlighted the substantial contribution of knowledge graph-derived features. Moreover, KG-CNNNTI was employed for virtual screening of natural products against Alzheimer's disease, resulting in 40 candidate compounds. 5 were supported by literature evidence, among which 3 were further validated *in vitro* assays.

1. Introduction

Drug development is a time-consuming and costly process, typically taking over 12 years and costing approximately \$2.6 billion from target identification to final market approval^{1,2}. Unfortunately, around 90% of drug candidates failed in clinical trials ultimately³. The identification of drug-target interactions (DTIs) plays a crucial role in various stages of drug development, including optimization of lead compounds⁴, drug repurposing⁵, and elucidation of drug side effects⁶. However, traditional DTI detection techniques are expensive and time-intensive. To address these limitations, computational approaches for prediction of DTIs have been developed as effective alternatives to experimental methods, offering new strategies for identification of potential drug candidates and their corresponding targets⁷⁻⁹.

Computational methods for DTI prediction can be broadly classified into three categories: structure-based, ligand-based, and hybrid approaches. Structure-based methods, such as mo-

lecular docking and molecular dynamics simulations, have been widely used for predicting drug-target binding affinity and in virtual screening¹⁰. However, their predictive capacity is significantly limited when the three-dimensional (3D) structure of the target protein is unavailable. Ligand-based approaches rely on the availability of active ligands for known targets. When there is insufficient bioactivity data for specific targets, these methods tend to suffer from limited generalization ability and predictive performance¹¹. Hybrid methods, which integrate both structural information of proteins and bioactivity data of ligands, have become increasingly important in DTI prediction. These approaches are typically categorized into two major types: network-based and artificial intelligence (AI)-based strategies. By leveraging both molecular structure and large-scale biological knowledge, these hybrid approaches help overcome the limitations of traditional structure- and ligand-based methods, particularly in cases with sparse or incomplete data. Furthermore, with the rapid accumulation of biomedical data and the continuous advancement of AI algorithms, network- and AI-driven hybrid methods are emerging as mainstream strategies in the field of computational drug discovery, offering new potential for improving the accuracy and generalizability of DTI prediction.¹²

Network-based methods construct large-scale DTI networks

* Corresponding author.

E-mail addresses: ruiwang@ecust.edu.cn (R. Wang); ytang234@ecust.edu.cn (Y. Tang)

(Y. Tang)

^A These authors contributed equally to this work.

and infer potential drug-target associations by propagating information through the network using various algorithms. Methods such as NBI¹³, SDTNBI¹⁴, bSDTNBI¹⁵, and wSDTNBI¹⁶ apply network inference techniques for DTI prediction. Although these approaches often demonstrate competitive performance, their predictive accuracy heavily depends on the completeness and quality of existing DTI data. HiSIF-DTA¹⁷, HGRL-DTA¹⁸ and MSF-DTA¹⁹ take biological network information as external prior knowledge, reducing the reliance on the DTI dataset to a certain extent. In contrast, AI-based methods involve feature representation of drugs and targets, followed by feature fusion and classification using machine learning or deep learning models. Traditional machine learning techniques, including Random forest (RF)²⁰, Multilayer Perceptron (MLP)²¹, and LightGBM²² have been widely applied in DTI prediction. In recent years, end-to-end deep learning architectures have demonstrated superior learning capabilities by automatically learning task-specific representations from molecular and sequence data, offering new strategies for DTI prediction. For instance, GraphDTA represents drugs as molecular graphs and proteins as one-hot encoded sequences, employing Graph Neural Networks (GNNs) for interaction prediction²³. TransformerCPI utilizes the Transformer-based architecture to learn high-dimensional representations of both protein sequences and compounds, leveraging self-attention mechanisms to capture long-range dependencies²⁴. DeepConv-DTI integrates Convolutional Neural Networks (CNNs) for local feature extraction and deep neural network model for modeling, thereby enabling the modeling complex DTI patterns more effectively²⁵.

Although current DTI prediction methods have achieved notable progress, several challenges still remain. One critical issue lies in the limited representation of proteins. Most existing approaches rely solely on one-dimensional amino acid sequences, neglecting 3D structural information and relevant biological functions. In fact, protein structure and interaction patterns fundamentally determine their unique biological functions, making the integration of functional information crucial for DTI prediction²⁶. Knowledge graph (KG) integrates diverse heterogeneous biomedical data, such as protein functions, molecular interactions, and disease associations, to comprehensively depict complex biomolecular relationships. By applying knowledge graph embedding (KGE) techniques, KGs can map entities and relations into a unified low-dimensional vector space, which compensate for the limitations of traditional models in representing protein functions. Recently, KG-based deep learning methods have made some progress^{12, 27}. However, most existing models rely heavily on known DTI data and often underutilize other biologic-

al information such as signaling pathways and disease phenotypes, leading to high sparsity and limited generalization capability.

To address the insufficient of protein functional representation, in this study we proposed a novel DTI prediction model, KG-CNNNTI (Fig. 1), aiming at enhancing the generalizability and computational efficiency of DTI prediction and providing a more robust tool for drug discovery. Specifically, this study extracted topological features of proteins from the knowledge graph utilizing the Node2Vec²⁸ algorithm and integrated them with sequence-based representations generated by ProteinBERT²⁹ to form comprehensive protein embeddings. For compound representation, we compared five types of molecular fingerprints along with a pre-trained representation method, Uni-Mol³⁰, which contains plenty of structural information. These features were then used as inputs to three machine learning models and one deep learning model for DTI prediction. The optimal model, KG-CNNNTI, was selected from 26 models based on 6 evaluation metrics, namely Accuracy, Precision, Recall, F1-Score, AUC and AUPR. The results on Davis and Metz dataset^{31, 32} demonstrated that KG-CNNNTI outperformed existing methods across several key metrics, including Precision, Recall, F1-Score and AUPR. Ablation studies confirmed the significant contribution of KG-derived features to prediction performance. Leveraging this capability, KG-CNNNTI successfully identified 40 potential multi-target compounds for Alzheimer's disease (AD). Among these, 5 compounds were validated by literature reports, and 3 ones were confirmed through *in vitro* experiments, demonstrating the model's promising screening performance and value of drug discovery.

2. Materials and methods

2.1. Data collection and preparation

2.1.1. Training dataset

The DTI data were collected from our previous research¹⁶, with filtering criteria including human (*Homo sapiens*) targets and activity types limited to K_i or K_d values. After removing duplicates and conflicting entries, a total of 282 342 high-quality DTI records were retained. RDKit was employed to preprocess the molecular structures, including desalting, charge neutralization, and standardization of SMILES representations, and molecules with invalid SMILES were excluded. The target proteins

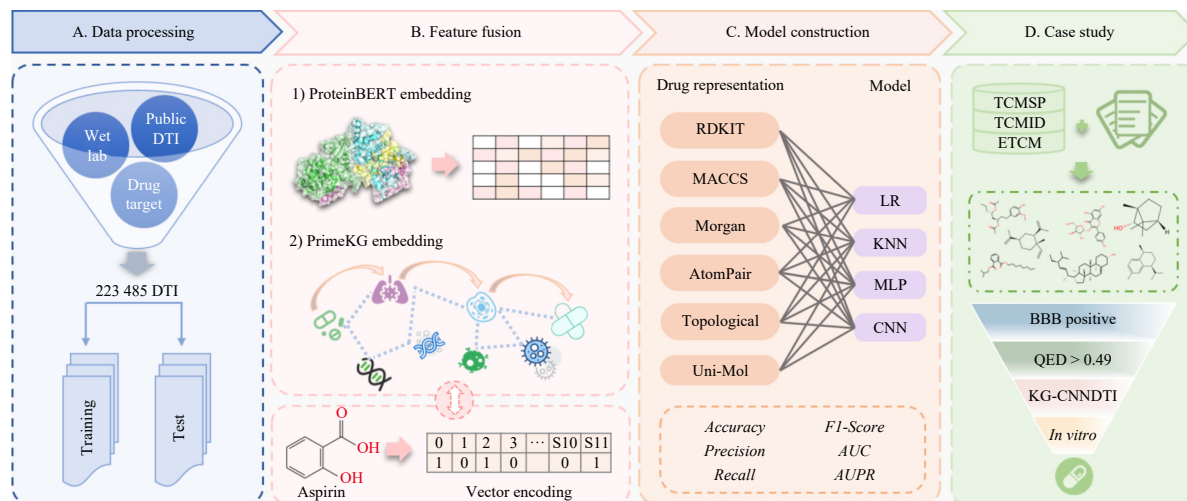


Fig. 1 The workflow of the KG-CNNNTI. (A) Data processing. (B) Feature fusion. (C) Model Construction. (D) Case study in anti-AD drug discovery from traditional Chinese medicine.

were then matched against nodes in the knowledge graph, and proteins not present in the graph were discarded, resulting in 223 485 DTI pairs used for model construction. Based on the data distribution provided with Supplementary Table S1, interactions with activity values $\leq 100 \text{ nmol}\cdot\text{L}^{-1}$ were labeled as positive and those $\geq 1000 \text{ nmol}\cdot\text{L}^{-1}$ as negative, resulting in a relatively balanced dataset with a positive-to-negative ratio of approximately 1 : 1.6. To avoid potential label noise caused by hard thresholding, interactions with activity values between 100 and 1000 $\text{nmol}\cdot\text{L}^{-1}$ were excluded. In order to assess the generalization capability of the model, the Davis dataset³¹ and Metz dataset³² were selected as external validation sets. Davis dataset contains 64 unique drugs and 379 unique protein targets, comprising 9 011 positive and 1 428 negative samples. Metz dataset contains 1 415 unique drugs, 135 unique proteins, 14 996 negative samples, and 16 160 positive samples.

2.1.2. Traditional Chinese medicine for anti-Alzheimer's disease

To explore the model's applicability in real-world drug discovery, we applied it to screen potential active compounds against potential targets of AD. A total of 4 282 AD-related publications were retrieved from CNKI and PubMed. Based on pharmacological and clinical evidence reported in the literature, traditional Chinese medicine (TCM), including herbal formulas, single herbs, and their active ingredients were identified after deduplication. The ingredients of the TCM were collected from the TCM-SP³³, ETCM³⁴, and TCMID³⁵ databases. To eliminate redundancy, compound names were first used to retrieve their corresponding SMILES structures from the PubChem database³⁶, and duplicates were then removed based on standardized SMILES using RDKit. In addition, to evaluate central nervous system drug-likeness, the blood-brain barrier (BBB) permeability of each compound was predicted using admetSAR 3.0³⁷, retaining only those predicted as BBB-positive. Furthermore, the quantitative estimation of drug-likeness (QED) was calculated using RDKit, and compounds with QED values > 0.49 were selected according to literature standards³⁸. Ultimately, 3 131 candidate compounds with potential activity were obtained.

2.2. Model architecture

2.2.1. Data representation

In DTI prediction task, the representation of drug molecules plays a pivotal role in determining model performance. To identify the most suitable molecular features, we explored six types of molecular representations, including five classical molecular fingerprints, namely MACCS, RDKit, AtomPair, Topological, and Morgan Fingerprint, as well as a pre-trained model Uni-Mol.

The classical molecular fingerprints were generated using RDKit, with 166-bit for MACCS and 512-bit for the remaining types. These fingerprints are binary vectors, where each bit represents the presence (1) or absence (0) of a specific molecular feature, such as a predefined substructure, atom path, or topological pattern. In contrast, Uni-Mol is a Transformer-based pre-trained molecular model that learns 512-dimensional continuous feature vectors by incorporating 3D atomic coordinates and geometric relationships. Unlike handcrafted fingerprints, Uni-Mol captures richer intra-molecular spatial information in a data-driven manner, offering improved generalizability in downstream tasks. By integrating and comparing both handcrafted and learned molecular features, this study provides insights into the impact of molecular representation on DTI prediction performance.

To obtain more comprehensive protein representations, we combined sequence-based features extracted by ProteinBERT with topological and semantic features derived from the biomed-

ical knowledge graph PrimeKG using Node2Vec. ProteinBERT effectively captures sequence dependencies but lacks contextual biological knowledge, while PrimeKG provides a large-scale, multi-level integration of protein-related biological entities. By fusing the complementary strengths of these two sources, we aimed to construct more informative and biologically meaningful protein embeddings. To formalize this integration, firstly, sequence level features were extracted using the pre-trained ProteinBERT model, which encodes each protein sequence into a 1024-dimensional dense vector $\mathbf{f}_{\text{BERT}} \in R^{1024}$. While it may contain redundant features and increase the computational burden on the model, principal component analysis (PCA) was applied to reduce the dimension into 128, resulting in $\mathbf{f}_{\text{PCA}} \in R^{128}$. Secondly, Node2vec was applied on PrimeKG to obtain a 128-dimensional network topology feature vector $\mathbf{f}_{\text{N2V}} \in R^{128}$. Finally, the complete protein representation, $\mathbf{f}_{\text{protein}}$, was constructed by concatenating the two vectors:

$$\mathbf{f}_{\text{protein}} = [\mathbf{f}_{\text{PCA}}; \mathbf{f}_{\text{N2V}}] \quad (1)$$

The protein feature matrix was concatenated with the molecular fingerprint matrix to form the complete DTI representation.

2.2.2. Model algorithm

In this study, we systematically compared the performance of conventional machine learning algorithms and deep learning models for DTI prediction. The machine learning models included Logistic Regression (LR), K-Nearest Neighbor (KNN), and MLP, while the deep learning model involved CNN.

Among the machine learning models, LR is a classical linear classifier that applies the sigmoid function to a linear combination of input features to estimate class probabilities. KNN is a non-parametric, instance-based method that classifies samples by majority vote from their k-nearest neighbors in the feature space. MLP is a typical feedforward neural network composed of input, hidden, and output layers, using ReLU activation to learn complex patterns and trained *via* backpropagation. All machine learning models were implemented using scikit-learn (version 1.2.2).

For the deep learning model, we developed a one-dimensional convolutional neural network using the PyTorch (version 2.0.1) framework to jointly learn representations from molecular and protein features for classification tasks. The input consisted of a 768-dimensional vector including 512 molecular features and 256 protein features. The network architecture included two convolution-pooling units for hierarchical feature extraction, followed by a fully connected layer for high-level integration and prediction. Model training was performed with the Adam optimizer (learning rate = 0.001) and binary cross-entropy loss, using a batch size of 32. An early stopping strategy was applied to reduce overfitting. To achieve optimal performance, we conducted a grid search over key hyperparameters, including out channels, kernel sizes, pooling sizes, hidden layer dimensions, and dropout rates.

The detailed search range of the model parameters were provided with Supplementary Table S2.

2.3. Model construction

After data processing, the dataset was divided into a training set and a test set at a ratio of 8 : 2 by scikit-learn. The training set was used for model training, while the test set was reserved for performance evaluation. To enhance the model's stability and generalization ability, 10-fold cross-validation was applied on the training set. Additionally, grid search was employed to optimize hyperparameters and identify the optimal parameter configuration. Given that this study addresses a binary classification task, it is crucial to select appropriate evaluation metrics to assess both

model accuracy and generalization capability. Each prediction result can be categorized into one of four outcomes: True Positive (TP), False Positive (FP), True Negative (TN), or False Negative (FN). Based on these outcomes, a series of statistical metrics were calculated to evaluate the model's performance from multiple perspectives.

Accuracy, defined as the proportion of correctly classified samples, is suitable for datasets with balanced class distributions.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (2)$$

Precision refers to the proportion of correctly predicted positive samples among all samples predicted as positive, reflecting the reliability of the positive predictions.

$$Precision = \frac{TP}{TP + FP} \quad (3)$$

Recall measures the proportion of true positive samples correctly identified by the model, indicating the model's ability to capture positive instances.

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

F1-score is utilized to provide a balanced assessment, which harmonically combines Precision and Recall.

$$F1-Score = \frac{2 * Precision * Recall}{Precision + Recall} \quad (5)$$

The Receiver Operating Characteristic (ROC) curve plots the true positive rate against the false positive rate across different threshold settings. The area under the ROC curve (AUC) ranges from 0 to 1, with higher values indicating better classification performance. In contrast, the area under the Precision-Recall curve (AUPR) is more informative in cases of class imbalance, offering a more accurate assessment of the model's ability to identify positive samples when their prevalence is low. An AUPR closer to 1 indicates better performance in recognizing positive instances.

To further assess the effectiveness of the proposed KG-CNNDTI framework, we conducted experiments on the Davis dataset and Metz Dataset and compared the results with those of existing state-of-the-art models. In addition, ablation studies were conducted on protein features to evaluate the contributions of different protein representation strategies.

2.4. Baseline models

To further evaluate the performance of the proposed model, three representative baseline methods were selected for comparison: MCANet³⁹, MolTrans⁴⁰, HyperAttentionDTI⁴¹, and BEACON⁴². These baseline models adopt distinct strategies for DTI prediction. Specifically, MCANet employs multi-scale convolutional modules to extract features from drug and protein sequences, and incorporates a cross-attention mechanism to strengthen the modeling of DTI. MolTrans utilizes a Transformer-based architecture to encode SMILES representations and protein sequences, effectively capturing long-range dependencies between semantic substructures. HyperAttentionDTI constructs a hypergraph to model complex and heterogeneous molecular interactions, and leverages multi-head attention to aggregate information across different relation types. BEACON uses GNN and CNN models to obtain structural features for compounds and proteins, respectively, and also utilizes a knowledge graph to acquire their knowledge-based features. A bilinear attention network is then constructed to fuse these extracted features. These baseline models reflect diverse modeling strategies in DTI prediction and serve as strong comparative references for assessing the effectiveness of the proposed approach.

2.5. Model validation

2.5.1. Molecular docking

The 3D structures of proteins were obtained from the RCSB Protein Data Bank (PDB)⁴³ and comprehensive details were presented in Supplementary Table S3. The molecular docking procedure was performed using Schrödinger software package (version 2021, Schrödinger, LLC) with the following configuration settings. At first, protein structures were preprocessed using the Protein Preparation Wizard module which included structure correction, hydrogen addition, removal of water molecules, and energy minimization. Water molecules located more than 5 Å away from the ligand were removed. Ionization states and tautomers were generated using the Epik module within a pH range of 5–9. Protein structures were optimized using the OPLS_2005 force field. Secondly, small-molecule ligands were prepared using the LigPrep module, including hydrogen addition, salt removal, generation of three-dimensional structures, tautomers, and all possible stereoisomers. Thirdly, the docking grid box was generated utilizing the Receptor Grid Generation module, with the co-crystallized ligand set as the center of the docking region. Finally, molecular docking was conducted using the Glide module in Standard Precision (SP) mode, employing the OPLS_2005 force field for the docking process.

2.5.2. Experimental validation

AD is pathologically characterized by extracellular deposition of Aβ plaques and intracellular neurofibrillary tangles. Increasing evidence suggests that Aβ deposition is closely associated with neuroinflammation⁴⁴. Aβ can bind to receptors such as Toll-like receptors (TLRs) and the receptor for advanced glycation end-products (RAGE), thereby activating surrounding microglia⁴⁵. Activated microglia and astrocytes can phagocytose Aβ and exert neuroprotective effects, however, failure to effectively clear Aβ aggregates leads to sustained inflammatory responses and the release of various pro-inflammatory mediators⁴⁶. In this study, a neuroinflammatory model induced by LPS was established to evaluate the therapeutic potential of candidate compounds for AD.

BV-2 mouse microglial were cultured in high-glucose DMEM supplemented with 10% fetal bovine serum and 1% penicillin/streptomycin at 37 °C in a humidified incubator containing 5% CO₂. To evaluate the effect of drug candidacies on BV-2 cell viability, cells were seeded in 96-well plates, after 24 hours, treated with serum-free medium (control) or drug candidacies diluted in serum-free medium (experimental), followed by 24 h incubation. For the LPS-induced model, BV-2 cells were pretreated with drug candidacies for 2 h, then exposed to 1 μg·mL⁻¹ Lipopolysaccharide (LPS) in serum-free medium, and incubated for another 24 h; the control group received equal volumes of serum-free medium. MTT assay and Nitric oxide (NO) production assay carried out to assess cell viability and anti-inflammatory effect respectively. Detailed experimental procedures are provided in Supplementary Information 1.

3. Results

3.1. Data distribution

After thorough data cleaning and deduplication processes, a total of 223 485 DTI samples were retained for subsequent model development and evaluation. These samples including 1 173 unique protein targets and 84 856 distinct small-molecule compounds. Among them, 84 535 were labeled as positive interactions, while 138 950 were categorized as negative samples. The

detailed statistics and distribution of the dataset were summarized in Table 1.

Table 1 The distribution of dataset

Type	Positive	Negative	Total
K_i	77 800	112 624	190 424
K_d	6 735	26 326	33 061
Total	84 535	138 950	223 485

To assess the consistency and representativeness of the dataset distribution, we applied Principal Component Analysis (PCA) to reduce the dimensionality of the compound feature space to

two dimensions. This analysis was conducted DTI data across six different molecular representation schemes: AtomPair, MACCS, Morgan, RDKit, Topological fingerprints, and the Uni-Mol pre-trained embeddings. The resulting 2D visualizations are illustrated in Fig. 2. As shown in the figure, the training and test sets exhibit substantial overlap across all molecular representation types, indicating that the samples are drawn from similar distributions. This high degree of distributional consistency suggests that the random split strategy adopted during dataset construction is reasonable and unbiased. Furthermore, the overlapping patterns between the training and test sets enhance the credibility of performance evaluations, ensuring that the test set serves as a valid proxy for assessing the model's generalization ability in real-world scenarios.

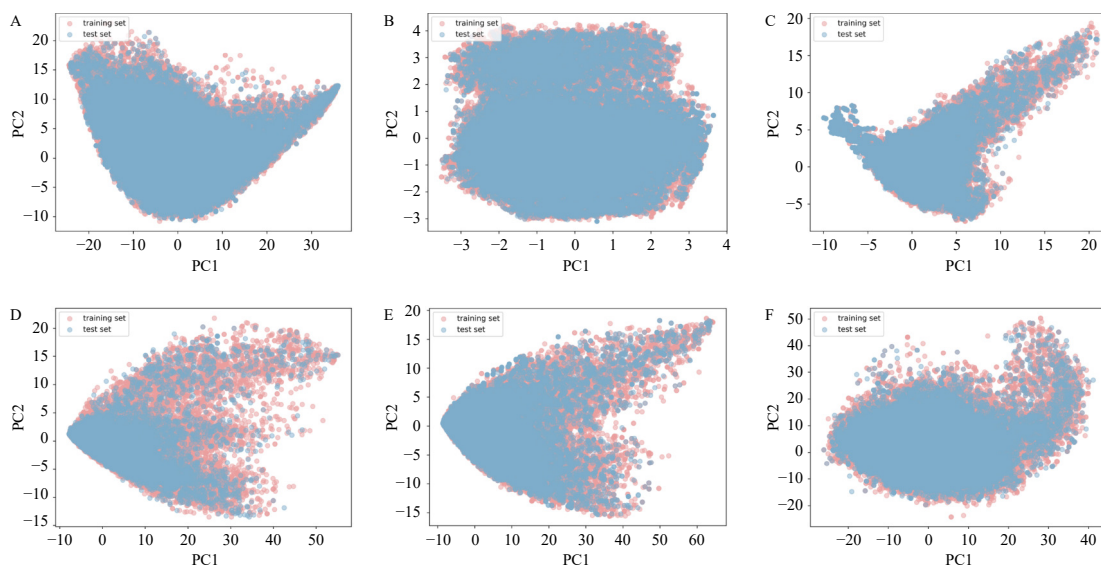


Fig. 2 The DTI data distribution across six different molecular representations. (A) AtomPair fingerprint. (B) MACCS fingerprint. (C) Morgan fingerprint. (D) RDKit fingerprint. (E) Topological fingerprint. (F) Uni-Mol pre-training representation.

3.2. Model performance analysis

3.2.1. The performance on the test set

In this study, a total of 24 DTI prediction models were constructed by combining six molecular representation methods with four prediction algorithms. *AUC* is an important metric for evaluating a classification model's ranking ability, generalization performance, and overall effectiveness. Model performance on the test set was evaluated using the *AUC* value, as illustrated in Fig. 3. All models achieved *AUC* values above 0.8, indicating strong generalization ability and predictive performance. Among all the models, those based on Morgan fingerprints consistently

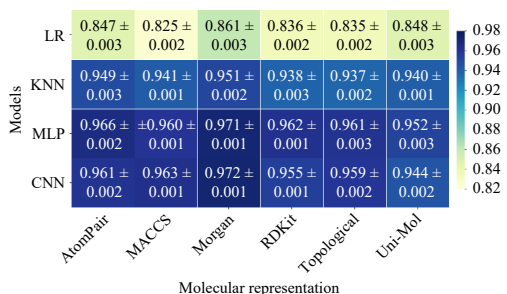


Fig. 3 The *AUC* values of different models on the test set. All the *AUC* values were more than 0.80, which indicates that the models have good performance. The Morgan fingerprint was the best of all molecular representations in our DTI dataset.

outperformed others across all four algorithm types. Morgan fingerprints encode the topological environment of each atom and its neighbors within a defined radius into fixed-length bit vectors, effectively capturing both local and global structural features of molecules⁴⁷.

Table 2 presented the top 10 models ranked by *AUC* on the test set. Among them, the model that integrated Morgan fingerprints with a CNN achieved the best overall performance across multiple evaluation metrics. Specifically, it reached an Accuracy of 0.921, a Recall of 0.931, an *AUC* of 0.972, and an AUPR of 0.970, ranking first among all models. Its precision was 0.914, ranking second. Based on these results, this model was identified as the most effective binary classification model for DTI prediction in this study and hereafter referred to as KG-CNNDTI. We conducted a sensitivity analysis on the key parameters of KG-CNNDTI. The results indicate that the model maintains stable *AUC*, AUPR, and F1-Score under parameter variations, demonstrating its strong robustness. Detailed results are provided in Supplementary Information 1. In subsequent work, KG-CNNDTI would be compared with existing state-of-the-art models to further validate its effectiveness and advantages in DTI prediction tasks.

3.2.2. Ablation experiment of protein representation

To further evaluate the role and contribution of knowledge graph embeddings in DTI prediction tasks, ablation studies were conducted based on the best-performing models among six different molecular representation strategies. The selected models

Table 2 The TOP 10 models with AUC values in the test set.

Model	Accuracy	Precision	Recall	F1-Score	AUC	AUPR
CNN_Morgan	0.921 (± 0.001)	0.914 (± 0.002)	0.931 (± 0.001)	0.917 (± 0.003)	0.972 (± 0.001)	0.970 (± 0.001)
MLP_Morgan	0.917 (± 0.004)	0.920 (± 0.002)	0.913 (± 0.001)	0.921 (± 0.005)	0.971 (± 0.001)	0.967 (± 0.002)
MLP_AtomPair	0.911 (± 0.001)	0.906 (± 0.003)	0.919 (± 0.001)	0.912 (± 0.003)	0.966 (± 0.002)	0.962 (± 0.004)
CNN_MACCS	0.906 (± 0.003)	0.895 (± 0.001)	0.920 (± 0.004)	0.911 (± 0.002)	0.963 (± 0.001)	0.959 (± 0.003)
MLP_RDKit	0.910 (± 0.001)	0.904 (± 0.003)	0.919 (± 0.004)	0.912 (± 0.002)	0.962 (± 0.001)	0.955 (± 0.001)
MLP_Topological	0.910 (± 0.002)	0.911 (± 0.001)	0.909 (± 0.003)	0.916 (± 0.002)	0.961 (± 0.003)	0.955 (± 0.002)
CNN_AtomPair	0.907 (± 0.001)	0.905 (± 0.001)	0.909 (± 0.003)	0.907 (± 0.003)	0.961 (± 0.002)	0.953 (± 0.003)
MLP_MACCS	0.905 (± 0.001)	0.914 (± 0.003)	0.896 (± 0.001)	0.905 (± 0.003)	0.960 (± 0.001)	0.956 (± 0.003)
CNN_Topological	0.902 (± 0.001)	0.890 (± 0.001)	0.917 (± 0.005)	0.903 (± 0.002)	0.959 (± 0.002)	0.952 (± 0.002)
CNN_RDKit	0.897 (± 0.001)	0.884 (± 0.001)	0.914 (± 0.004)	0.899 (± 0.001)	0.955 (± 0.001)	0.949 (± 0.001)

included four based on MLP and two based on CNN. In these experiments, the knowledge graph embedding vectors and ProteinBERT-derived protein sequence features were separately removed from the input to assess the performance differences with and without each feature type. The results of the ablation experiments are shown in Fig. 4. It was observed that removing the knowledge graph embeddings led to a noticeable decline in model performance across all test cases, indicating the significant contribution of these features. Notably, the performance degradation caused by removing the knowledge graph features was consistently greater than that observed when excluding protein sequence embeddings. Specifically, compared to the removal of protein sequence features, the exclusion of knowledge graph features resulted in performance drops averagely by folds of 6.3, 3.4, 7.3, 5.0, 3.6, and 3.1 in terms of Accuracy, Precision, Recall, F1-Score, AUC, and AUPR, respectively. This finding further highlights the value of knowledge graph embeddings in capturing the latent semantic associations between drugs and targets. In summary, knowledge graph-based features not only enhance the model's ability to represent complex biomolecular relationships

but also, in many cases, outperform traditional sequence-based features. They serve as a critical factor in improving the performance of DTI prediction models.

To evaluate the effectiveness of PrimeKG for protein representation, we compared it with Hetionet, a widely used biomedical knowledge graph. Specifically, we extracted protein node embeddings from Hetionet using the same Node2Vec-based embedding strategy and replaced the PrimeKG-based embeddings in our model. The results, presented in Fig. 4, show that the performance of the model using Hetionet embeddings was consistently lower across multiple evaluation metrics, compared to the model using PrimeKG-based embeddings. This suggests that the broader and more diverse biological relationships captured in PrimeKG provide richer contextual information for protein function modeling. These findings support our decision to adopt PrimeKG as the primary knowledge graph in this study.

3.2.3. Comparison with existing models

The comparative results with existing models, as illustrated in Fig. 5, demonstrated that KG-CNNDTI consistently outper-

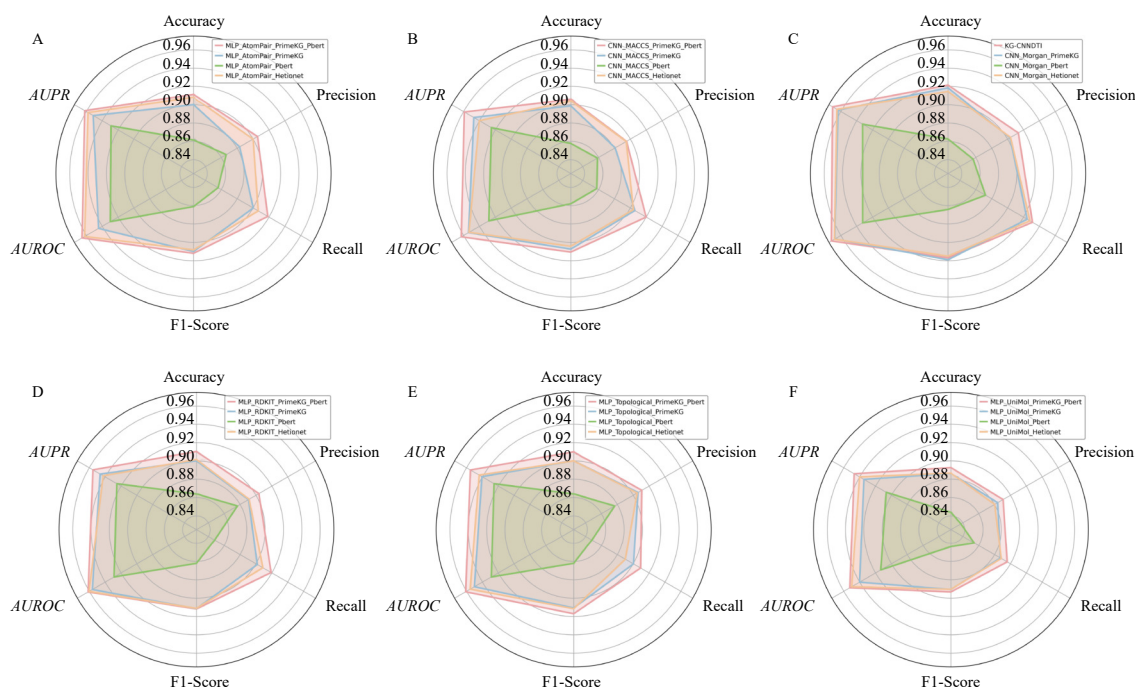


Fig. 4 The ablation experiments on models with different molecular representation. (A) AtomPair fingerprint. (B) MACCS fingerprint. (C) Morgan fingerprint. (D) RDKit fingerprint. (E) Topological fingerprint. (F) Uni-Mol pre-training representation. The name of the model indicates the architecture of the model, and the different modules of the model are separated by “_”. The first one is the prediction algorithm, including CNN and MLP, the second one is the molecular representation method, including 5 types of molecular fingerprints and a pre-trained model, and the third one is the protein representation method, including knowledge graph features integration with protein sequence features (PrimeKG_Pbert), knowledge graph features (PrimeKG or Hetionet) only, and protein sequence features (Pbert) only.

formed all baseline models in terms of Precision, Recall, F1-Score and AUPR, highlighting its superior predictive capability on the validation datasets. In particular, the high Precision indicates that the model produces fewer false positives, which is crucial in DTI prediction tasks where the number of negative pairs far exceeds the positives. Minimizing false positives helps reduce the cost of downstream experimental validation. Similarly, the superior AUPR reflects the model's ability to maintain high precision across various recall levels. A high F1-score indicates that the model has robust performance. Although the performance of KG-CNNNTI on Davis dataset, specifically in terms of accuracy and AUC is slightly lower than that of the top-performing models, the overall results validate the effectiveness and robustness of the proposed framework. These findings suggest that integrating knowledge graph embeddings and sequence-based features provides valuable biological context and improves the practical utility of KG-CNNNTI in drug discovery and screening applications.

3.3. Prediction of anti-AD ingredients

To demonstrate the practical utility of the proposed KG-CNNNTI model, we conducted a case study focusing on multi-target drug discovery for AD. Increasing evidence suggests that the complex and multifactorial nature of AD could not be effectively addressed by single-target therapies, which often exhibit limited clinical efficacy⁴⁸. Consequently, there is a growing need to develop multi-target drug of modulating several key pathological pathways simultaneously. Compared to traditional drug-target prediction models, KG-CNNNTI incorporates a knowledge graph enriched with protein functional information, making it more competitive in screening multi-target and multi-pharmacological activity drugs for Alzheimer's disease. In line with this rationale, we curated 13 well-known AD-associated targets based on prior studies⁴⁹, representing a broad spectrum of molecular mechanisms implicated in AD progression. To identify potential multi-

target compounds, we collected and preprocessed a total of 3 131 structurally unique compounds from the databases and literature. The preprocessing pipeline included structural deduplication and SMILES standardization. To assess the chemical diversity of the dataset, we calculated the Tanimoto similarity coefficients between all compound pairs using MACCS fingerprints. The analysis yielded a low average Tanimoto similarity score of 0.157 (Fig. 6), reflecting high structural heterogeneity within the dataset which is a desirable characteristic for virtual screening.

Subsequently, we applied the KG-CNNNTI model to systematically predict potential DTI between the TCM compounds and the 13 AD-related targets. The predicted interaction probabilities exhibited a bimodal distribution, with values clustering around [0, 0.1] and [0.9, 1] (Supplementary Fig. S2), indicating the model's robust binary classification capability. By applying a high-confidence threshold (probability ≥ 0.9), we identified 40 candidate compounds predicted to interact with all 13 targets. These compounds may exert therapeutic effects via synergistic modulation of multiple AD-related pathways, thereby offering promising leads for the development of multi-target therapeutics against Alzheimer's disease.

Further analysis indicated that 5 of the 40 candidate compounds are closely associated with major pathological pathways of AD, such as amyloid- β aggregation, neuroinflammatory responses, oxidative stress, and cognitive impairment (Fig. 7). Based on compound availability and experimental feasibility, three were ultimately selected for *in vitro* validation. Specifically, talatisamine (Compound 8), an alkaloid from *Aconitum*, is known for blocking potassium channels and alleviating A β -induced cytotoxicity⁵⁰. Molecular docking revealed that its binding energies with four AD-related targets were below $-5 \text{ kcal}\cdot\text{mol}^{-1}$, indicating strong binding affinity. In the LPS-induced inflammation model, Talatisamine reduced NO production by 19.06% compared to the model group (Figs. 8A, 8D). Karakoline (Compound 9), a nicotinic receptor agonist, has been reported to improve cognitive function in AD mouse models⁵¹. It exhibited favorable binding with

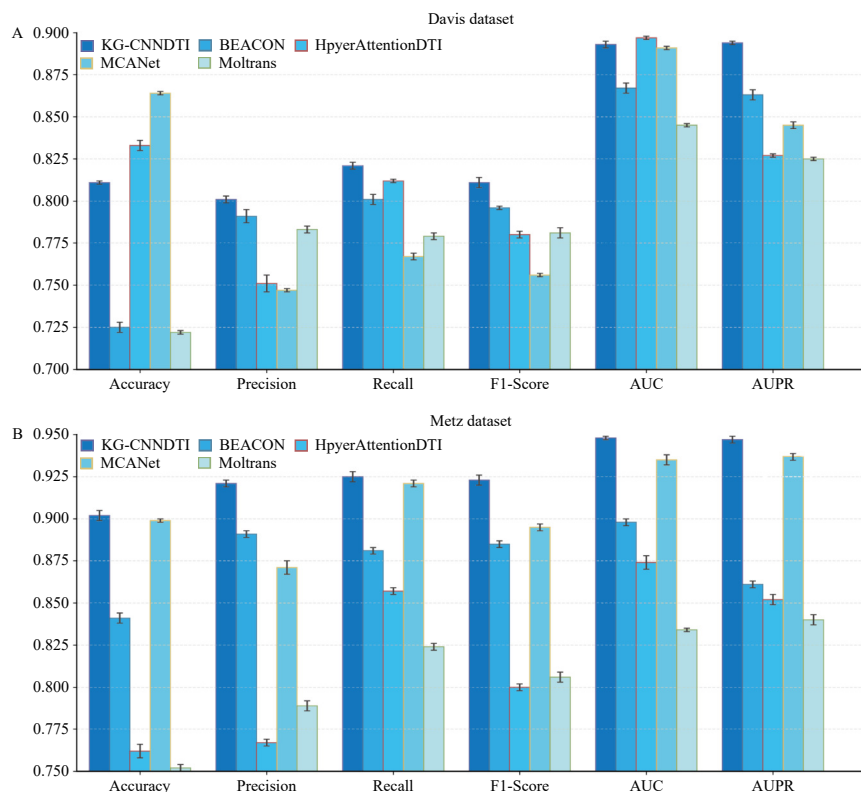


Fig. 5 KG-CNNNTI outperforms four baseline models on the Davis and Metz datasets. KG-CNNNTI achieves superior performance on Precision, Recall, F1-Score, and AUPR on Davis (A), and excels across six evaluation metrics (exceeding 0.9) on Metz (B).

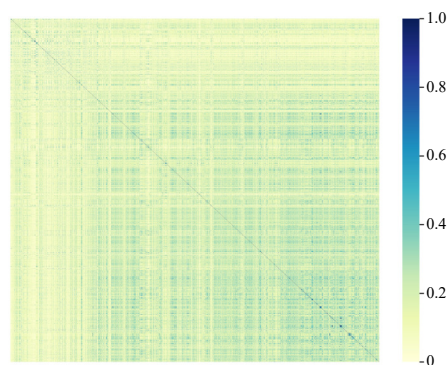


Fig. 6 The structural similarity analysis of 3131 ingredients of TCM. The average Tanimoto similarity of the 3131 compounds was 0.157, indicating molecular diversity.

eight targets and reduced NO production by 11.15% in co-culture experiments (Figs. 8C, 8F). Hapepunine (Compound **18**), a steroidal alkaloid, displayed potent anti-inflammatory activity with an IC_{50} value of $20.85 \mu\text{mol}\cdot\text{L}^{-1}$ in inhibiting LPS-induced NO release⁵² and showed strong binding with seven targets in molecular docking. Gitogenin (Compound **19**) activates the AMPK signaling pathway to induce apoptosis and autophagy⁵³. In cell experiments, it reduced NO production by 10.64% relative to the model group (Figs. 8B, 8E). Withanolide B (Compound **38**) has been reported to reduce A β aggregation and mitigate neurotoxicity⁵⁴, and demonstrated strong binding with five targets in docking simulations. It showed docking scores below $-5 \text{ kcal}\cdot\text{mol}^{-1}$ with eight AD-related targets. The whole results of molecular docking were provided with Supplementary Table S4 and all the molecular information were provided with Supplementary Fig. S3. In general, this study successfully identified a set of promising natural products with well-defined pharmacological mechanisms through virtual screening. These results highlight the practical value and robustness of the KG-CNNNTI model and provide novel leads for AD drug discovery.

4. Discussion

DTI prediction plays a crucial role in lead optimization, drug repositioning, and elucidating mechanisms of adverse drug reactions. However, current studies often exhibit limitations in representing protein functions. To address this issue, we proposed the KG-CNNNTI framework, which incorporates protein knowledge graph embeddings enriched with biological semantics to effectively complement and enhance protein representations. This

approach not only improves the accuracy and applicability of DTI prediction but also enhances its generalization performance in practical drug screening tasks. Compared to other existing methods, our KG-CNNNTI has the following innovation.

Firstly, the model substantially improves protein representation by integrating structural and functional knowledge from heterogeneous sources. Unlike conventional models that rely solely on protein sequences, KG-CNNNTI combines topological features extracted from biomedical knowledge graphs *via* Node2Vec with contextualized embeddings from ProteinBERT. Ablation studies confirmed that the fusion of knowledge graph-derived embeddings led to a significant boost in predictive performance, highlighting their ability to capture biologically meaningful relationships.

Secondly, the model demonstrates superior performance across multiple evaluation metrics on the Davis and Metz dataset when compared with state-of-the-art baselines. Specifically, KG-CNNNTI achieves higher performance in the metrics of Precision, Recall, F1-Score and AUPR, indicating its ability to accurately identify true interactions while minimizing false positives. Furthermore, the strong performance across different metrics indicates that KG-CNNNTI can reliably identify meaningful drug-target interactions, highlighting its robustness and generalizability across diverse datasets.

Finally, the applicability of KG-CNNNTI was further validated in a case study on AD. The model successfully identified several candidate compounds with potential therapeutic relevance. Notably, Talatisamine (Molecule 8) was identified anti-AD by alleviating A β -induced cytotoxicity. In *in vitro* experimental validation, Talatisamine exhibited significant anti-inflammatory activity, with a 19.06% reduction in NO production compared to the model group. This suggests its possible role in alleviating AD-related neuroinflammation. These findings not only support the model's generalization capability but also demonstrate its potential to facilitate multi-target drug discovery for complex diseases.

Despite the promising performance of KG-CNNNTI, several limitations remain. First, the model lacks interpretability, as its deep learning structure functions largely as a "black-box", limiting its usefulness in explaining predictions. Although knowledge graphs provide biological context, there is no explicit mechanism to interpret the learned embeddings. Second, the model architecture can be further optimized. It does not yet fully leverage advanced feature interaction techniques such as matrix factorization^{55,56} or attention mechanisms^{57,58}, which could enhance semantic alignment and performance on complex datasets. Third, the quality and completeness of the underlying knowledge graph are potential bottlenecks. Limitations in node diversity, edge at-

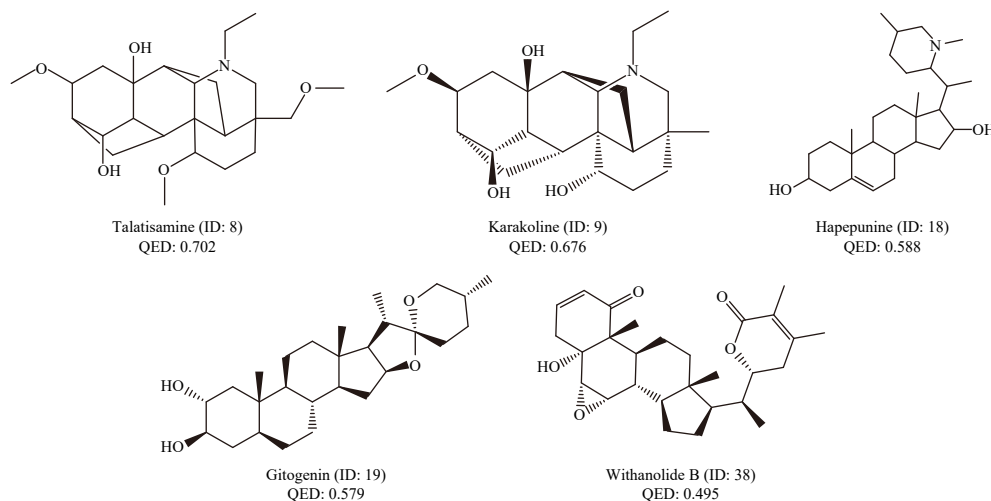


Fig. 7 5 molecules with potential anti-AD activity. The QED values for each molecule were calculated using RDKit.

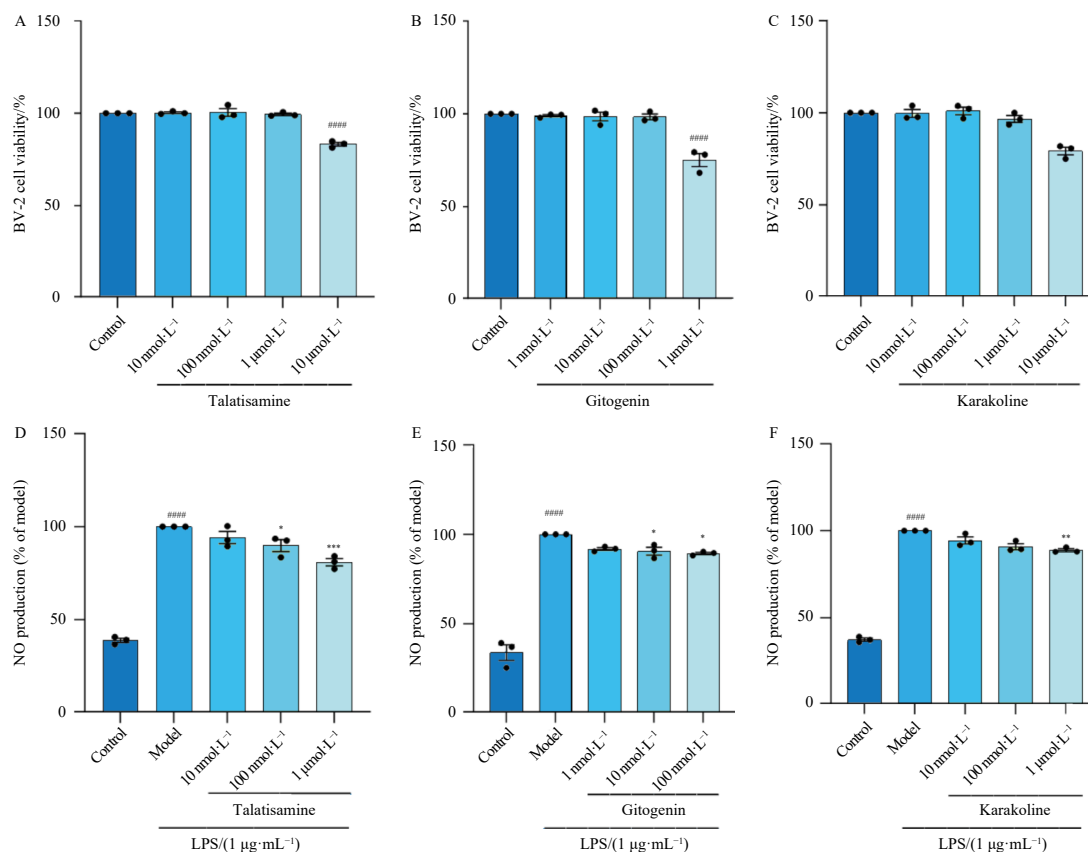


Fig. 8 Results of cell experiments. (A-C) Cytotoxicity assays showed that talatisamine and karakoline significantly reduced BV-2 cell viability at 10 $\mu\text{mol}\cdot\text{L}^{-1}$, while gitogenin showed significant cytotoxicity at 1 $\mu\text{mol}\cdot\text{L}^{-1}$. Non-toxic concentrations were used in subsequent experiments. (D-F) Talatisamine (1 $\mu\text{mol}\cdot\text{L}^{-1}$), gitogenin (100 $\text{nmol}\cdot\text{L}^{-1}$), and karakoline (1 $\mu\text{mol}\cdot\text{L}^{-1}$) significantly inhibited LPS-induced NO production, with inhibition rates of 19.06%, 10.64%, and 11.15% compared with model group, respectively. All data represent three independent experiments. Statistical analysis was performed using GraphPad Prism 10.1.2. Brown-Forsythe test assessed variance homogeneity, Dunnett's multiple comparison test was used to compare data of more than two groups, and the *t*-test was used to compare two sets of data. (**** $P < 0.0001$ vs Control group; * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$ vs Model group).

tributes, and update frequency may constrain representation capacity. Future improvements may involve integrating attention-based interpretability, self-supervised learning⁵⁹, graph convolution models^{60,61}, and dynamic biomedical knowledge graphs⁶² to enhance model transparency, robustness, and adaptability. Finally, Although the LPS-induced inflammation model provided initial evidence for the anti-inflammatory effects of the compounds, it does not fully reflect the complex pathology of AD. Key mechanisms such as A β aggregation, tau hyperphosphorylation, and neurodegeneration were not evaluated. Future studies will include additional assays, such as A β aggregation and neuroprotection tests, to further validate their anti-AD potential.

5. Conclusions

In this study, we proposed a knowledge graph-enhanced framework named KG-CNNDTI for DTI prediction. By integrating protein sequence features with embedding features derived from biomedical knowledge graphs, the model constructed a more comprehensive and biologically relevant representation of protein targets. In addition, six types of molecular representations and four modeling algorithms were systematically explored to evaluate their performance in DTI prediction tasks. The best-performing model, KG-CNNDTI, achieved an *AUC* value of 0.972 on the test set, demonstrating excellent predictive capability. Ablation studies indicated that knowledge graph embedding features contributed more significantly to prediction performance than sequence-based features alone. Furthermore, KG-CNNDTI outperformed several state-of-the-art methods across multiple evaluation metrics on different datasets, highlighting its robustness and generalizability. In practical drug screening, KG-CNNDTI

identified 40 potential multi-target compounds. Notably, several of these compounds have been validated in the previous literature and our own *in vitro* experiments, further confirming the reliability and practical utility of the model. Overall, this study demonstrated the value of integrating knowledge graphs with molecular and sequence-based features for DTI prediction, and would provide a promising strategy for identification of active compounds from traditional Chinese medicine in the context of neurodegenerative disease drug discovery.

Funding

This work was supported by the National Natural Science Foundation of China (Nos. 82173746 and U23A20530) and Shanghai Frontiers Science Center of Optogenetic Techniques for Cell Metabolism (Shanghai Municipal Education Commission).

Supporting information

The supporting information includes a more detailed introduction to the training data and results, as well as a detailed description of the experimental reagents. This information can be obtained from the author *via* email.

Declaration of competing interest

These authors have no conflict of interest to declare.

References

- 1 New drug costs soar to \$2.6 billion. *Nat Biotech.* 2014;32(12):1176.

- 2 Zhang K, Yang X, Wang Y, et al. Artificial intelligence in drug development. *Nat Med*. 2025;31(1):45-59. <https://doi.org/10.1038/s41591-024-03434-4>.
- 3 Takebe T, Imai R, Ono S. The current status of drug discovery and development as originated in United States academia: the influence of industrial and academic collaboration on drug discovery and development. *Clin Transl Sci*. 2018;11(6):597-606. <https://doi.org/10.1111/cts.12577>.
- 4 Tang B, He F, Liu D, et al. AI-directed design of novel targeted covalent inhibitors against SARS-CoV-2. *Biomolecules*. 2022;12(6):746. <https://doi.org/10.3390/biom12060746>.
- 5 Ajmal A, Mahmood A, Hayat C, et al. Computer-assisted drug repurposing for thymidylate kinase drug target in monkeypox virus. *Front Cell Infect Microbiol*. 2023;13:1159389. <https://doi.org/10.3389/fcimb.2023.1159389>.
- 6 Yu Z, Wu Z, Zhou M, et al. mtADENet: a novel interpretable method integrating multiple types of network-based inference approaches for prediction of adverse drug events. *Comput Biol Med*. 2024;168:107831. <https://doi.org/10.1016/j.combiomed.2023.107831>.
- 7 Xu L, Ru X, Song R. Application of machine learning for drug-target interaction prediction. *Front Genet*. 2021;12:680117. <https://doi.org/10.3389/fgene.2021.680117>.
- 8 Dhakal A, McKay C, Tanner JJ, et al. Artificial intelligence in the prediction of protein-ligand interactions: recent advances and future directions. *Brief Bioinform*. 2022;23(1):bbab476. <https://doi.org/10.1093/bib/bbab476>.
- 9 Liao Q, Zhang Y, Chu Y, et al. Application of artificial intelligence in drug-target interactions prediction: a review. *npj Biomed Innov*. 2025;2(1):1. <https://doi.org/10.1038/s44385-024-00003-9>.
- 10 Xu Z, Guan L, Wang Y, et al. Discovery of a novel PLK1 inhibitor with high inhibitory potency using a combined virtual screening strategy. *J Enzyme Inhib Med Chem*. 2025;40(1):2467798. <https://doi.org/10.1080/14756366.2025.2467798>.
- 11 Sharma V, Wakode S, Kumar H. Chapter 2-structure-and ligand-based drug design: concepts, approaches, and challenges. *Chemoinformatics and Bioinformatics in the Pharmaceutical Sciences*. Amsterdam: Academic Press, 2021: 27-53.
- 12 Ye Q, Hsieh CY, Yang Z, et al. A unified drug-target interaction prediction framework based on knowledge graph and recommendation system. *Nat Commun*. 2021;12(1):6775. <https://doi.org/10.1038/s41467-021-27137-3>.
- 13 Cheng F, Liu C, Jiang J, et al. Prediction of drug-target interactions and drug repositioning via network-based inference. *PLoS Comput Biol*. 2012;8(5):e1002503. <https://doi.org/10.1371/journal.pcbi.1002503>.
- 14 Wu Z, Cheng F, Li J, et al. SDTNBI: an integrated network and chemoinformatics tool for systematic prediction of drug-target interactions and drug repositioning. *Brief Bioinform*. 2017;18(2):333-347. <https://doi.org/10.1093/bib/bbw012>.
- 15 Wu Z, Lu W, Wu D, et al. In silico prediction of chemical mechanism of action via an improved network-based inference method. *Brit J Pharmacol*. 2016; 173(23):3372-3385. <https://doi.org/10.1111/bph.13629>.
- 16 Wu Z, Ma H, Liu Z, et al. wSDTNBI: a novel network-based inference method for virtual screening. *Chem Sci*. 2022; 13(4): 1060-1079.
- 17 Bi X, Zhang S, Ma W, et al. HISIF-DTA: a hierarchical semantic information fusion framework for drug-target affinity prediction. *IEEE J Biomed Health Inform*. 2025;29(3):1579-1590. <https://doi.org/10.1109/JBHI.2023.3334239>.
- 18 Chu Z, Huang F, Fu H, et al. Hierarchical graph representation learning for the prediction of drug-target binding affinity. *Inform Science*. 2022;613(C): 507-523. <https://doi.org/10.1016/j.ins.2022.09.043>.
- 19 Ma W, Zhang S, Li Z, et al. Predicting drug-target affinity by learning protein knowledge from biological networks. *IEEE J Biomed Health Inform*. 2023;27(4):2128-2137. <https://doi.org/10.1109/JBHI.2023.3240305>.
- 20 Shi H, Liu S, Chen J, et al. Predicting drug-target interactions using Lasso with random forest based on evolutionary information and chemical structure. *Genomics*. 2019;111(6):1839-1852. <https://doi.org/10.1016/j.ygeno.2018.12.007>.
- 21 Li F, Zhang Z, Guan J, et al. Effective drug-target interaction prediction with mutual interaction neural network. *Bioinformatics*. 2022;38(14):3582-3589. <https://doi.org/10.1093/bioinformatics/btac377>.
- 22 Peng Y, Zhao S, Zeng Z, et al. LGBMDF: a cascade forest framework with LightGBM for predicting drug-target interactions. *Front Microbiol*. 2022;13: 1092467. <https://doi.org/10.3389/fmicb.2022.1092467>.
- 23 Nguyen T, Le H, Quinn TP, et al. GraphDTA: predicting drug-target binding affinity with graph neural networks. *Bioinformatics*. 2021;37(8):1140-1147. <https://doi.org/10.1093/bioinformatics/btaa921>.
- 24 Chen L, Tan X, Wang D, et al. TransformerCPI: improving compound-protein interaction prediction by sequence-based deep learning with self-attention mechanism and label reversal experiments. *Bioinformatics*. 2020;36(16): 4406-4414. <https://doi.org/10.1093/bioinformatics/btaa524>.
- 25 Lee I, Keum J, Nam H. DeepConv-DTI: prediction of drug-target interactions via deep learning with convolution on protein sequences. *PLoS Comput Biol*. 2019;15(6):e1007129. <https://doi.org/10.1371/journal.pcbi.1007129>.
- 26 Sivley RM, Dou X, Meiler J, et al. Comprehensive analysis of constraint on the spatial distribution of missense variants in human protein structures. *Am J Hum Genet*. 2018;102(3):415-426. <https://doi.org/10.1016/j.ajhg.2018.01.017>.
- 27 Li N, Yang Z, Wang J, et al. Drug-target interaction prediction using knowledge graph embedding. *iScience*. 2024;27(6):109393. <https://doi.org/10.1016/j.isci.2024.109393>.
- 28 Grover A, Leskovec J. node2vec: scalable feature learning for networks. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2016: 855-864. <https://doi.org/10.1145/2939672.2939754>.
- 29 Brandes N, Ofer D, Peleg Y, et al. ProteinBERT: a universal deep-learning model of protein sequence and function. *Bioinformatics*. 2022;38(8):2102-2110. <https://doi.org/10.1093/bioinformatics/btac020>.
- 30 Zhou G, Gao Z, Ding Q, et al. Uni-mol: a universal 3d molecular representation learning framework. *International Conference on Learning Representations* 2023.
- 31 Davis MI, Hunt JP, Herrgard S, et al. Comprehensive analysis of kinase inhibitor selectivity. *Nat Biotechnol*. 2011;29(11):1046-1051. <https://doi.org/10.1038/nbt.1990>.
- 32 Metz JT, Johnson EF, Soni NB, et al. Navigating the kinome. *Nat Chem Biol*. 2011;7(4):200-202. <https://doi.org/10.1038/nchembio.530>.
- 33 Ru J, Li P, Wang J, et al. TCMSp: a database of systems pharmacology for drug discovery from herbal medicines. *J Cheminform*. 2014;6:13. <https://doi.org/10.1186/1758-2946-6-13>.
- 34 Xu HY, Zhang YQ, Liu ZM, et al. ETCM: an encyclopaedia of traditional Chinese medicine. *Nucleic Acids Res*. 2019;47(D1):D976-D982. <https://doi.org/10.1093/nar/gky987>.
- 35 Huang L, Xie D, Yu Y, et al. TCMID 2.0: a comprehensive resource for TCM. *Nucleic Acids Res*. 2018;46(D1):D1117-D1120. <https://doi.org/10.1093/nar/gkx1028>.
- 36 Kim S, Chen J, Cheng T, et al. PubChem 2025 update. *Nucleic Acids Res*. 2024;53(D1):D1516-D1525. <https://doi.org/10.1093/nar/gkae1059>.
- 37 Gu Y, Yu Z, Wang Y, et al. admetSAR3.0: a comprehensive platform for exploration, prediction and optimization of chemical ADMET properties. *Nucleic Acids Res*. 2024;52(W1):W432-W438. <https://doi.org/10.1093/nar/gkae298>.
- 38 Bickerton GR, Paolini GV, Besnard J, et al. Quantifying the chemical beauty of drugs. *Nat Chem*. 2012;4(2):90-98. <https://doi.org/10.1038/nchem.1243>.
- 39 Brians J, Zhang X, Zhang X, et al. MCANet: shared-weight-based MultiheadCrossAttention network for drug-target interaction prediction. *Brief Bioinform*. 2023;24(2):bbad082. <https://doi.org/10.1093/bib/bbad082>.
- 40 Huang K, Xiao C, Glass LM, et al. MolTrans: molecular Interaction transformer for drug-target interaction prediction. *Bioinformatics*. 2021;37(6):830-836. <https://doi.org/10.1093/bioinformatics/btaa880>.
- 41 Zhao Q, Zhao H, Zheng K, et al. HyperAttentionDTI: improving drug-protein interaction prediction by sequence-based deep learning with attention mechanism. *Bioinformatics*. 2022;38(3):655-662. <https://doi.org/10.1093/bioinformatics/btab715>.
- 42 Tao W, Lin X, Liu Y, et al. Bridging chemical structure and conceptual knowledge enables accurate prediction of compound-protein interaction. *BMC Biol*. 2024;22(1):248. <https://doi.org/10.1186/s12915-024-02049-y>.
- 43 Burley SK, Bhikadiya C, Bi C, et al. RCSB protein data bank: powerful new tools for exploring 3D structures of biological macromolecules for basic and applied research and education in fundamental biology, biomedicine, biotechnology, bioengineering and energy sciences. *Nucleic Acids Res*. 2020; 49(D1):D437-D451. <https://doi.org/10.1093/nar/gkaa1038>.
- 44 Calsolaro V, Edison P. Neuroinflammation in Alzheimer's disease: current evidence and future directions. *Alzheimers Dement*. 2016;12(6):719-732. <https://doi.org/10.1016/j.jalz.2016.02.010>.
- 45 Heneka MT, Kummer MP, Latz E. Innate immune activation in neurodegenerative disease. *Nat Rev Immunol*. 2014;14(7):463-477. <https://doi.org/10.1038/nri3705>.
- 46 Onyango IG, Jauregui GV, Čarná M, et al. Neuroinflammation in Alzheimer's disease. *Biomedicines*. 2021;9(5):524. <https://doi.org/10.3390/biomedicines9050524>.
- 47 Rogers D, Hahn M. Extended-connectivity fingerprints. *J Chem Inf Model*. 2010;50(5):742-754. <https://doi.org/10.1021/ci100050t>.
- 48 Zheng Q, Wang X. Alzheimer's disease: insights into pathology, molecular mechanisms, and therapy. *Protein Cell*. 2025;16(2):83-120. <https://doi.org/10.1093/procel/pwae026>.
- 49 Yue C, Chen B, Pan F, et al. TCnet: a novel strategy to predict target combination of Alzheimer's disease via network-based methods. *J Chem Inf Model*. 2025;65(7):3866-3878. <https://doi.org/10.1021/acs.jcim.5c00172>.
- 50 Yang X, Xin Y, Gu Y, et al. Total alkaloids of *Aconitum carmichaelii* Debx alleviate cisplatin-induced acute renal injury by inhibiting inflammation and oxidative stress related to gut microbiota metabolism. *Phytomedicine*. 2024;135:156128. <https://doi.org/10.1016/j.phymed.2024.156128>.
- 51 Nie H, Wang Z, Zhao W, et al. New nicotinic analogue ZY-1 enhances cognitive functions in a transgenic mice model of Alzheimer's disease. *Neurosci Lett*. 2013;537:29-34. <https://doi.org/10.1016/j.neulet.2013.01.001>.
- 52 Wang L, Jiang Y, Yaseen A, et al. Steroidal alkaloids from the bulbs of *Fritillaria pallidiflora* Schrenk and their anti-inflammatory activity. *Bioorg Chem*. 2021;112:104845. <https://doi.org/10.1016/j.bioorg.2021.104845>.
- 53 Liu T, Li Y, Sun J, et al. Gtogenin suppresses lung cancer progression by inducing apoptosis and autophagy initiation through the activation of AMPK signaling. *Int Immunopharmacol*. 2022;111:108806. <https://doi.org/10.1016/j.intimp.2022.108806>.
- 54 Balkrishna A, Bhattacharya K, Shukla S, et al. Neuroprotection by polyherbal medicine divya-medha-vati against scopolamine-induced cognitive impairment through modulation of oxidative stress, acetylcholine activity, and cell signaling. *Mol Neurobiol*. 2024;61(3):1363-1382. <https://doi.org/10.1007/s12035-023-03601-7>.
- 55 Liu T, Wang S, Zhang Y, et al. TIWMFLP: two-tier interactive weighted matrix factorization and label propagation based on similarity matrix fusion for drug-disease association prediction. *J Chem Inf Model*. 2024;64(22):8641-8654. <https://doi.org/10.1021/acs.jcim.4c01589>.
- 56 Wang S, Liu T, Ren C, et al. Predicting potential small molecule-miRNA associations utilizing truncated Schatten p-norm. *Brief Bioinform*. 2023; 24(4):bbad234. <https://doi.org/10.1093/bib/bbad234>.
- 57 Zhang L, Wang CC, Zhang Y, et al. GPCNDTA: prediction of drug-target binding affinity through cross-attention networks augmented with graph features and pharmacophores. *Comput Biol Med*. 2023;166:107512. <https://doi.org/10.1016/j.combiomed.2023.107512>.
- 58 Li J, Bi X, Ma W, et al. MHAN-DTA: a multiscale hybrid attention network for drug-target affinity prediction. *IEEE J Biomed Health Inform*. 2024;2024:1-21. <https://doi.org/10.1109/JBHI.2024.3518619>.
- 59 Liu Z, Chen Q, Lan W, et al. SSLDTI: a novel method for drug-target interaction prediction based on self-supervised learning. *Artif Intell Med*. 2024;149:102778. <https://doi.org/10.1016/j.artmed.2024.102778>.
- 60 Peng W, Chen T, Liu H, et al. Improving drug response prediction based on two-space graph convolution. *Comput Biol Med*. 2023;158:106859. <https://doi.org/10.1016/j.combiomed.2023.106859>.
- 61 Lan W, Zhou G, Chen Q, et al. Contrastive clustering learning for multi-behavior recommendation. *ACM Trans Inf Syst*. 2024; 43(1): Article 18. <https://doi.org/10.1145/3698192>.
- 62 Lan W, Tang Z, Liu M, et al. The large language models on biomedical data analysis: a survey. *IEEE J Biomed Health Inform*. 2025;29(6):4486-4497. <https://doi.org/10.1109/JBHI.2025.3530794>.