FusionDTI: Fine-grained Binding Discovery with Token-level Fusion for Drug-Target Interaction

Anonymous ACL submission

Abstract

Predicting drug-target interaction (DTI) is critical in the drug discovery process. Despite remarkable advances in recent DTI models through the integration of representations from diverse drug and target encoders, such models often struggle to capture the fine-grained interactions between drugs and protein, i.e. the binding of specific drug atoms (or substructures) and key amino acids of proteins, which is crucial for understanding the binding mechanisms and optimising drug design. To address this issue, this paper introduces a novel model, called FusionDTI, which uses a token-level Fusion module to effectively learn fine-grained information for Drug-Target Interaction. In particular, our FusionDTI model uses the SELF-IES representation of drugs to mitigate sequence fragment invalidation and incorporates the structure-aware (SA) vocabulary of target proteins to address the limitation of amino acid sequences in structural information, additionally leveraging pre-trained language models extensively trained on large-scale biomedical datasets as encoders to capture the complex information of drugs and targets. Experiments on three well-known benchmark datasets show that our proposed FusionDTI model achieves the best performance in DTI prediction compared with eight existing state-of-the-art baselines. Furthermore, our case study indicates that FusionDTI could highlight the potential binding sites, enhancing the explainability of the DTI prediction.¹

1 Introduction

011

012

017

022

040

041

The task of predicting drug-target interactions (DTI) plays a pivotal role in the drug discovery progress, as it helps identify potential therapeutic effects of drugs on biological targets facilitating the development of effective treatments (Askr et al., 2023). DTI fundamentally relies on the binding of specific drug atoms (or substructures) and key amino acids of proteins (Schenone et al., 2013). In particular, each binding site is an interaction between a single amino acid and a single drug atom, which we refer to as a fine-grained interaction. For instance, Figure 1 B demonstrates the interaction between HIV-1 protease and the drug lopinavir. A critical component of this interaction is the formation of a hydrogen bond between a ketone group in lopinavir (represented in the SELFIES (Krenn et al., 2022) notation as [C][=O]) and the side chain of an aspartate residue Asp25 (i.e. Dd) within the protease (Brik and Wong, 2003; Chandwani and Shuter, 2008). Therefore, capturing such finegrained interaction information during the fusion of drug and target representations is crucial for building effective DTI prediction models (Wu et al., 2022; Peng et al., 2024; Zeng et al., 2024).

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

To obtain representations of drugs and targets for the DTI task, some previous studies (Lee et al., 2019; Nguyen et al., 2021) have used graph neural networks (GNNs) or convolutional neural networks (CNNs) using a fixed-size window, potentially leading to a loss of contextual information, especially when drugs and targets are in a longterm sequence. These models directly concatenate the representations together to make predictions without considering fine-grained interactions. More recently, some computational models (Huang et al., 2021; Bai et al., 2023) employed the fusion module (e.g. Deep Interactive Inference Network (DIIN) (Gong et al., 2018) and Bilinear Attention Network (BAN) (Kim et al., 2018)) to obtain finegrained interaction information and the 3-mer approach that binds three amino acids together as a target binding site to address the lack of structural information in the amino acid sequence. While useful for highlighting possible regions of interaction, these models do not offer the sufficient granularity needed to gauge the specifics of binding sites, as each binding site only contains one

¹The complete code and datasets are available in the software section of the submission.



Figure 1: **A**. Illustration of the FusionDTI model: frozen encoder, fusion module and classifier. The token-level fusion (TF) focuses on fine-grained interactions between tokens within and across sequences. **B**. This is a token-level interaction instance of HIV-1 protease and lopinavir. Lopinavir forms a hydrogen bond with residue Dd (Asp25) in the active site of the protease via its ketone molecule ([C][=O]). **C**. The attention map of TF visualises the weight between tokens, indicating the contribution of each drug atom and residue to the final prediction result.

residue (Schenone et al., 2013). Therefore, obtaining contextual representations of drugs and targets and capturing fine-grained interaction information for DTI remains challenging.

086

091

100

101

102

103

104

107

109

110

111

112

113

114

115

To address these challenges, we propose a novel model (called FusionDTI) with a Token-level Fusion (TF) module for an effective learning of fine-grained interactions between drugs and targets. In particular, our FusionDTI model utilises two pre-trained language models (PLMs), namely Saport (Su et al., 2023) as the protein encoder that is able to integrate both residue tokens with structure token; and SELFormer (Yüksel et al., 2023) as the drug encoder to ensure that each drug is valid and contains structural information. To effectively learn fine-grained information from these contextual representations of drugs and targets, we explore two strategies for the TF module, i.e. Bilinear Attention Network (BAN) (Kim et al., 2018) and Cross Attention Network (CAN) (Li et al., 2021; Vaswani et al., 2017), to find the best approach for integrating the rich contextual embeddings derived from Saport and SELFormer. We conduct a comprehensive performance comparison against eight existing state-of-the-art DTI prediction models. The results show that our proposed model achieves about 6% accuracy improvement over the best baseline on the BindingDB dataset. The main contributions of our study are as follows:

> • We propose FusionDTI, a novel model that leverages PLMs to encode drug SELFIES, as well as protein residues and structures for rich semantic representations and uses the token

level fusion to capture fine-grained interaction between drugs and targets effectively.

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

138

139

140

141

142

143

144

145

146

- We compare two TF modules: CAN and BAN and analyse the influence of fusion scales based on FusionDTI, demonstrating that CAN is superior for DTI prediction both in terms of effectiveness and efficiency.
- We conduct a case study of three drug-target pairs by FusionDTI to evaluate whether potential binding sites would be highlighted for the DTI prediction explainability.

2 Related Work

2.1 Drug and Protein Representation

For drug molecules, most existing methods represent the input by the Simplified Molecular Input Line Entry System (SMILES) (Weininger, 1988; Weininger et al., 1989). However, SMILES suffers from numerous problems in terms of validity and robustness, and some valuable information about the drug structure may be lost which may prevent the model from efficiently mining the knowledge hidden in the data (Krenn et al., 2022). To address the limitations of SMILES, we apply SELFIES, a string-based representation that circumvents the issue of robustness and that always generates valid molecular graphs for each character.

Regarding proteins, the conventional approach uses amino acid sequences as model inputs (Huang et al., 2021; Bai et al., 2023), overlooking the crucial structural information of the protein. Inspired by the SA vocabulary of SaProt (Su et al., 2023),

196

the SaProt enhances inputs by amalgamating each 147 residue of the amino acid sequence with a 3D geo-148 metric feature that is obtained by encoding protein 149 structure information using Foldseek (Van Kem-150 pen et al., 2024). This innovative combination 151 offers richer protein representations through the 152 SA vocabulary, contributing to the discovery of 153 fine-grained interactions. 154

2.2 Molecular and Protein Language Models

156

159

160

161

162

163

164

165

167

168

169

171

172

173

174

175

176

177

178

179

180

Molecular language models trained on the large-scale molecular corpus capture the subtleties of chemical structures and their biological activities, setting new standards in the encoding of chemical compounds achieving meaningful representations (Ying et al., 2021; Rong et al., 2020).
For example, MoLFormer (Ross et al., 2022) focused on leveraging the self-attention mechanism to interpret the complex, non-linear interactions within molecules, while SELFormer (Yüksel et al., 2023) employed SELFIES, ensuring valid and interpretable chemical structures.

Protein language models have revolutionized the way we understand and represent protein sequences, learning intricate patterns and features that define the protein functionality and interactions. ProtBERT (Elnaggar et al., 2021) and ESM (Lin et al., 2023) applied a transformer architecture to protein sequences, capturing the complex relationships between amino acids. Saport (Su et al., 2023, 2024) further enhanced this approach by integrating SA vocabularies to provide protein structure information.

3 Methodology

3.1 Model Architecture

Given a sequence-based input drug-target pair, the DTI prediction task aims to predict an interaction 182 probability score $p \in [0, 1]$ between the given drugtarget pair, which is typically achieved through 184 learning a joint representation F space from the 185 given sequence-based inputs. To address the DTI task and effectively capture fine-grained interac-187 tion, we proposed a novel model, called FusionDTI, which is a bi-encoder model (Liu et al., 2021) with a fusion module that fuses the representations of 190 191 drugs and targets. The overall framework of FusionDTI is illustrated in Figure 1 A. In general, 192 FusionDTI takes sequence-based inputs of drugs 193 and targets, which are encoded into token-level representation vectors by two frozen encoders. Then, 195

a fusion module fuses the representations to capture fine-grained binding information for a final prediction through a prediction head.

Input: The initial inputs of drugs and targets are string-based representations. For protein \mathcal{P} , the SA vocabulary (Su et al., 2023; Van Kempen et al., 2024) is employed, where each residue is replaced by one of 441 SA vocabularies that bind an amino acid to a 3D geometric feature to address the lack of structural information in amino acid sequences. For drug \mathcal{D} , as mentioned in the previous section, we use the SELFIES, which is a formal syntax that always generates valid molecular graphs (Krenn et al., 2022). We provide the steps and code to obtain SA and SELFIES in Appendix A.3.

Encoder: The proposed model contains two frozen encoders: Saport (Su et al., 2023) and SELF-ormer (Yüksel et al., 2023), which generate a drug representation **D** and a protein representation **P** separately. It is of note that FusionDTI is flexible enough to easily replace encoders with other PLMs or address SELFIES or SA representations that are unavailable. Furthermore, **D** and **P** are stored in memory for later-stage online training.

Fusion module: In developing FusionDTI, we have investigated two options for the fusion module: BAN and CAN to fuse representations, as indicated in Figure 2. The CAN is utilised to fuse each pair as D^* and P^* , and then concatenate them into one **F** for fine-grained binding information. For BAN, we need to obtain bilinear attention maps and generate **F** through the bilinear pooling layer.

Prediction head: Finally, we obtain the probability score p of the DTI prediction by a multilayer perceptron (MLP) classifier trained with the binary cross-entropy loss, i.e. $p = MLP(\mathbf{F})$.

Since the encoders and the fusion module constitute the key components of our FusionDTI model, we will describe them in detail in the following.

3.2 Drug and Protein Encoders

Employing sequences with detailed biological functions and structures is a critical step in exploring the fine-grained binding of drugs and targets. For drugs, SMILES is the most commonly used input sequence but suffers from invalid sequence segments and potential loss of structural information (Krenn et al., 2022). To address the limitations, we transform SMILES into SELFIES, a formal grammar that generates a valid molecular graph for each element (Krenn et al., 2022). Besides, to address the lack of structural information in the



Figure 2: **BAN:** In step 1, the bilinear attention map is obtained by a bilinear interaction modelling via transformation matrices. In step 2, the joint representation \mathbf{F} is generated using the attention map by bilinear pooling via the shared transformation matrices \mathbf{U} and \mathbf{V} . **CAN:** It fuses protein and drug representations through multi-head, self-attention and cross-attention. Then fused representations \mathbf{P}^* and \mathbf{D}^* are concatenated into \mathbf{F} after mean pooling.

amino acid sequences, we utilise the SA sequence of targets to combine each amino acid with an SA vocabulary by Foldseek (Van Kempen et al., 2024).

247

249

251

257

258

262

267

270

271

PLMs have shown promising achievements in the biomedical domain leveraging transformers since they pay attention to contextual information and are pre-trained on large-scale biomedical databases. Therefore, we utilise Saport (Su et al., 2023) as a protein encoder to encode protein input \mathcal{P} of both the SA sequence and amino acid sequence. Meanwhile, SELFormer (Yüksel et al., 2023) is used as our drug encoder to encode the drug SELFIES input \mathcal{D} . Then these encoded protein representation P and drug representation D are further used as inputs for the later fusion module (Subsection 3.3). These rich contextual representations ensure that we can explore the finegrained binding information effectively. To further justify this, we also compare our encoders with other existing protein language models (such as ESM-2 (Lin et al., 2023)) and molecular language models (such as MoLFormer (Ross et al., 2022) and ChemBERTa-2 (Ahmad et al., 2022)), and the results can be found in Appendix A.6.

3.3 Fusion Module

In order to capture the fine-grained binding information between a drug and a target, our FusionDTI model applies a fusion module to learn token-level interactions between the token representations of drugs and targets encoded by their respective encoders. As shown in Figure 2, two fusion modules are investigated to fuse representations: the Bilinear Attention Network (Kim et al., 2018) and the Cross Attention Network (Vaswani et al., 2017).

280

281

284

285

287

289

292

293

294

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

3.3.1 Bilinear Attention Network (BAN)

Motivated by DrugBAN (Bai et al., 2023), our model considers BAN (Kim et al., 2018) as an option to learn pairwise fine-grained interactions between drug $\mathbf{D} \in \mathbb{R}^{M \times \phi}$ and target $\mathbf{P} \in \mathbb{R}^{N \times \rho}$, denoted as FusionDTI-BAN. For BAN as indicated in Figure 2, bilinear attention maps are obtained by a bilinear interaction modelling to capture pairwise weights in step 1, and then the bilinear pooling layer to extract a joint representation **F**. The equation of BAN is shown below:

$$\mathbf{F} = \text{BAN}(\mathbf{P}, \mathbf{D}; Att)$$

= SumPool($\sigma(\mathbf{P}^{\top}\mathbf{U}) \cdot Att \cdot \sigma(\mathbf{D}^{\top}\mathbf{V}), s),$ (1)

where $\mathbf{U} \in \mathbb{R}^{N \times K}$ and $\mathbf{V} \in \mathbb{R}^{M \times K}$ are transformation matrices for representations. SumPool is an operation that performs a one-dimensional and non-overlapped sum pooling operation with stride *s* and $\sigma(\cdot)$ denotes a non-linear activation function with ReLU(·). *Att* $\in \mathbb{R}^{\rho \times \phi}$ represents the bilinear attention maps using the Hadamard product and matrix-matrix multiplication and is defined as:

$$Att = ((\mathbf{1} \cdot \mathbf{q}^{\top}) \circ \sigma(\mathbf{P}^{\top}\mathbf{U})) \cdot \sigma(\mathbf{V}^{\top}\mathbf{D}), \qquad (2)$$

Here, $\mathbf{1} \in \mathbb{R}^{\rho}$ is a fixed all-ones vector, $\mathbf{q} \in \mathbb{R}^{K}$ is a learnable weight vector and \circ denotes the Hadamard product. In this way, pairwise interactions contribute sub-structural pairs to predictions.

BAN captures the token-level interactions between the protein and drug representations without considering the relationships within each sequence itself, which may limit its ability to understand deeper contextual dependencies.

324

326

327

329

338

339

341 343 342

344

347

351

354

358

3.3.2 Cross Attention Network (CAN)

Inspired by ProST (Xu et al., 2023), we also con-312 sider CAN as our fusion module to learn fine-313 grained interaction information of drugs and targets. 314 We denote our FusionDTI model that uses a CAN fusion module as FusionDTI-CAN. By processing $\mathbf{D} \in \mathbb{R}^{m imes h}$ and $\mathbf{P} \in \mathbb{R}^{n imes h}$ separately, the fused 317 drug $\mathbf{D}^* \in \mathbb{R}^{m \times h}$ and target $\mathbf{P}^* \in \mathbb{R}^{n \times h}$ represen-318 tations are obtained. To synthesise the fine-grained 319 joint representation F, we employ a pooling aggregation strategy for both D^* and P^* independently 321 and then concatenate them as shown in Figure 2. The process is described by the following equation: 323

$$\mathbf{F} = \text{Concat}[\text{MeanPool}(\mathbf{D}^*), \text{MeanPool}(\mathbf{P}^*)], \quad (3)$$

where MeanPool calculates the element-wise mean of all tokens across the sequence dimension, and Concat denotes the concatenation of the resulting mean vectors. In this context, the multi-head, self-attention and cross-attention mechanisms are used to refine the representations of each residue and atom as below:

$$\mathbf{D}^{*} = \frac{1}{2} \left[MHA(\mathbf{Q}_{d}, \mathbf{K}_{d}, \mathbf{V}_{d}) + MHA(\mathbf{Q}_{p}, \mathbf{K}_{d}, \mathbf{V}_{d}) \right],$$
(4)

$$\mathbf{P}^{*} = \frac{1}{2} \left[MHA(\mathbf{Q}_{p}, \mathbf{K}_{p}, \mathbf{V}_{p}) + MHA(\mathbf{Q}_{d}, \mathbf{K}_{p}, \mathbf{V}_{p}) \right], \quad (5)$$

where $\mathbf{Q}_d, \mathbf{K}_d, \mathbf{V}_d \in \mathbb{R}^{m \times h}$ and $\mathbf{Q}_p, \mathbf{K}_p, \mathbf{V}_p \in \mathbb{R}^{n \times h}$ are the queries, keys and values for drug and target protein, respectively. And *MHA* denotes the Multi-head Attention mechanism. To guide this process, two distinct sets of projection matrices guide the attention mechanism as follows:

$$\mathbf{Q}_{d} = \mathbf{D}\mathbf{W}_{q}^{d}, \quad \mathbf{K}_{d} = \mathbf{D}\mathbf{W}_{k}^{d}, \quad \mathbf{V}_{d} = \mathbf{D}\mathbf{W}_{v}^{d}, \quad (6)$$
$$\mathbf{Q}_{p} = \mathbf{P}\mathbf{W}_{q}^{p}, \quad \mathbf{K}_{p} = \mathbf{P}\mathbf{W}_{k}^{p}, \quad \mathbf{V}_{p} = \mathbf{P}\mathbf{W}_{v}^{p}, \quad (7)$$

Here, the projection matrices $\mathbf{W}_q^d, \mathbf{W}_k^d, \mathbf{W}_v^d \in \mathbb{R}^{h \times h}$ and $\mathbf{W}_q^p, \mathbf{W}_k^p, \mathbf{W}_v^p \in \mathbb{R}^{h \times h}$ are used to derive the queries, keys and values, respectively.

In summary, our CAN module combines multihead, self-attention and cross-attention mechanisms to capture dependencies within individual sequences and between different sequences for a more nuanced understanding of interactions. In the results of Sections 4.3 and 4.5, we analyse and compare these two fusion strategies and different fusion scales in detail.

4 Experimental Setup and Results

4.1 Datasets and Baselines

Three public DTI datasets, namely BindingDB (Gilson et al., 2016), BioSNAP (Zitnik et al., 2018) and Human (Liu et al., 2015; Chen et al., 2020), are used for evaluation, where each dataset is split into training, validation, and test sets with a 7:1:2 ratio using two different splitting strategies: in-domain and cross-domain. For the in-domain split, the datasets are randomly divided. For the cross-domain setting, the datasets are split such that the drugs and targets in the test set do not overlap with those in the training set, making it a more challenging scenario where models must generalise to novel drug-target interactions. Since DTI is a binary classification task, we use AUROC (Bai et al., 2023; Huang et al., 2021) and AUPRC (Nguyen et al., 2021) as the major metrics to evaluate models' performance. In Appendix A.10, we report other evaluation metrics, including F1-score, Sensitivity, Specificity, and Matthews Correlation Coefficient (MCC) to provide a more comprehensive assessment.

We compare FusionDTI with eight baseline models in the DTI prediction task. These models include two traditional machine learning methods such as SVM (Cortes and Vapnik, 1995) and Random Forest (RF) (Ho, 1995), as well as five deep learning methods including DeepConv-DTI (Lee et al., 2019), GraphDTA (Nguyen et al., 2021), MolTrans (Huang et al., 2021), DrugBAN (Bai et al., 2023) and SiamDTI (Zhang et al., 2024). In addition, we also include the BioT5 (Pei et al., 2023) model, which is a biomedical pre-trained language model that could directly predict the DTI.

Furthermore, results on three additional benchmark datasets (DAVIS (Davis et al., 2011), KIBA (Tang et al., 2014), and DUD-E (Mysinger et al., 2012)) are reported, with comparisons to 8 task-specific baselines (Nga et al., 2025; Li et al., 2025). Further details regarding the datasets, baseline models, and the methodology for generating drug SELFIES and protein SA sequences are provided in Appendix A.3.

4.2 Evaluation of DTI Prediction

We start by comparing our FusionDTI model (FusionDTI-CAN and FusionDTI-BAN) with eight existing state-of-the-art baselines for DTI prediction on three widely used datasets. Table 1 reports the in-domain comparative results. In general, our FusionDTI-CAN model performs the best on all metrics across all three datasets. A key highlight from these results is the exceptional performance of FusionDTI-CAN on the BindingDB dataset, where FusionDTI-CAN demonstrates superior metrics

371

372

373

374

375

376

377

378

379

381

383

385

387

389

390

391

393

394

395

397

359

360

398 399 400

401 402

403

404

405

406

407

408

		BindingDB		Hur	nan		BioSNAP	
Method	AUROC	AUPRC	Accuracy	AUROC	AUPRC	AUROC	AUPRC	Accuracy
SVM	$.939 {\pm} .001$	$.928 {\pm} .002$	$.825 {\pm} .004$	$.940 {\pm} .006$	$.920 {\pm} .009$	$.862 {\pm} .007$	$.864 {\pm} .004$.777±.011
RF	$.942 {\pm} .011$	$.921 {\pm} .016$	$.880 {\pm} .012$	$.952 {\pm} .011$	$.953 {\pm} .010$	$.860 {\pm} .005$	$.886 {\pm} .005$	$.804 {\pm} .005$
DeepConv-DTI	$.945 {\pm} .002$	$.925 {\pm} .005$	$.882 {\pm} .007$	$.980 {\pm} .002$	$.981 {\pm} .002$	$.886 {\pm} .006$	$.890 {\pm} .006$	$.805 {\pm} .009$
GraphDTA	$.951 {\pm} .002$	$.934 {\pm} .002$	$.888 {\pm} .005$	$.981 {\pm} .001$	$.982 {\pm} .002$	$.887 {\pm} .008$	$.890 {\pm} .007$	$.800 {\pm} .007$
MolTrans	$.952 {\pm} .002$	$.936 {\pm} .001$	$.887 {\pm} .006$	$.980 {\pm} .002$	$.978 {\pm} .003$	$.895 {\pm} .004$	$.897 {\pm} .005$	$.825 {\pm} .010$
DrugBAN	$.960 {\pm} .001$	$.948 {\pm} .002$	$.904 {\pm} .004$	$.982 {\pm} .002$	$.980 {\pm} .003$	$.903 {\pm} .005$	$.902 {\pm} .004$	$.834 {\pm} .008$
SiamDTI	$.961 {\pm} .002$	$.945 {\pm} .002$	$.890 {\pm} .006$	$.970 {\pm} .002$	$.969 {\pm} .003$	$.912 {\pm} .005$	$.910 {\pm} .003$	$.855 {\pm} .004$
BioT5	$.963 {\pm} .001$	$.952 {\pm} .001$	$.907 {\pm} .003$	$\underline{.989 {\pm} .001}$	$\underline{.985 {\pm} .002}$	$\underline{.937 {\pm} .001}$	$\underline{.937 {\pm} .004}$	$\underline{.874} \pm .001$
FusionDTI-BAN	$.975 {\pm} .002$	$.976 {\pm} .002$	$.933 {\pm} .003$	$.984 {\pm} .002$	$.984 {\pm} .003$	$.923 {\pm} .002$	$.921 {\pm} .002$	$.856 {\pm} .001$
FusionDTI-CAN	.989±.002	.990±.002	.961±.002	$.991 {\pm} .002$	$\textbf{.989}{\pm}\textbf{.002}$	$.951 {\pm} .002$	$.952 {\pm} .002$	$\textbf{.889}{\pm}\textbf{.002}$

Table 1: In-domain performance comparison of FusionDTI and the baselines on the BindingDB, Human and BioSNAP datasets (**Best**, Second Best).

		BindingDB		Hui	nan		BioSNAP	
Method	AUROC	AUPRC	Accuracy	AUROC	AUPRC	AUROC	AUPRC	Accuracy
SVM	$.490 {\pm} .015$	$.460 {\pm} .001$	$.531 {\pm} .009$	$.621 \pm .036$	$.637 {\pm} .009$	$.602 {\pm} .005$	$.528 {\pm} .005$	$.513 {\pm} .011$
RF	$.493 {\pm} .021$	$.468 {\pm} .023$	$.535 {\pm}.012$	$.642 {\pm} .011$	$.663 {\pm} .050$	$.590 {\pm} .015$	$.568 {\pm} .018$	$.499 {\pm} .004$
GraphDTA	$.536 {\pm} .015$	$.496 {\pm} .029$	$.472 {\pm} .009$	$.822 {\pm} .009$	$.759 {\pm} .006$	$.618 {\pm} .005$	$.618 {\pm} .008$	$.535 {\pm} .024$
DeepConv-DTI	$.527 {\pm} .038$	$.499 {\pm} .035$	$.490 {\pm} .027$	$.761 {\pm} .016$	$.628 {\pm} .022$	$.645 {\pm} .022$	$.642 {\pm} .032$	$.558 {\pm} .025$
MolTrans	$.554 {\pm} .024$	$.511 {\pm} .025$	$.470 {\pm} .004$	$.810 {\pm} .021$	$.745 {\pm} .034$	$.621 \pm .015$	$.608 {\pm} .022$	$.546 {\pm} .032$
DrugBAN	$.604 {\pm} .027$	$.570 {\pm} .047$	$.509 {\pm} .021$	$.833 {\pm} .020$	$.760 {\pm} .031$	$.685 {\pm} .044$	$.713 {\pm} .041$	$.565 {\pm} .056$
SiamDTI	$.627 {\pm} .027$	$.571 {\pm} .024$	$.563 {\pm} .033$	$\textbf{.863}{\pm}\textbf{.019}$	$.807 \pm .040$	$.718 {\pm} .055$	$.725 {\pm} .054$	$.623 {\pm} .070$
BioT5	$.651 {\pm} .002$	$.653 {\pm} .003$	$.621 {\pm} .005$	$\underline{.856 {\pm}.003}$.853±.003	$.720{\pm}.008$	$.718 {\pm} .004$	$.715 {\pm} .009$
FusionDTI-BAN	$.659 {\pm} .002$	$.663 {\pm} .002$	$.633 {\pm} .003$	$.784 {\pm} .002$	$.790 {\pm} .003$	$.723 {\pm} .002$	$.721 {\pm} .002$	$.756 {\pm} .001$
FusionDTI-CAN	.681±.005	.680±.012	.652±.005	$.801 {\pm} .037$	$.803 {\pm} .032$.748±.021	.766±.017	.734±.012

Table 2: Cross-domain performance comparison of FusionDTI and the baselines on the BindingDB, Human and BioSNAP datasets (**Best**, <u>Second Best</u>).

across the board: an AUROC of 0.989, an AUPRC of 0.990, and an accuracy of 96.1%. Note that the main difference between the FusionDTI-CAN model and others is the fusion strategy. Furthermore, despite FusionDTI-BAN and DrugBAN both utilising the same BAN module, FusionDTI-BAN consistently outperforms DrugBAN on all datasets.

410 411

412

413

414

415

416

However, in-domain classification using random 417 splits holds limited practical significance. Thus, we 418 also evaluate the more challenging cross-domain 419 DTI prediction, where the training data and the 420 test data contain distinct drugs and targets. This 421 setting precludes the use of known drug or target 422 features when making predictions on the test data. 423 As shown in Table 2, the performance of all mod-424 els is diminished compared to the in-domain set-425 ting due to the reduced availability of information. 426 Nevertheless, the FusionDTI-CAN model demon-427 428 strates outstanding performance in cross-domain DTI prediction on the BindingDB and BioSNAP 429 datasets, highlighting its robustness in predicting 430 novel drug-target interactions. For instance, on the 431 BindingDB dataset, FusionDTI-CAN achieves the 432

highest metrics with an AUROC of 0.675 and an AUPRC of 0.676. This underscores the effectiveness of the model's fusion strategy in diverse and challenging scenarios. Similarly, despite sharing the BAN module, FusionDTI-BAN continues to outperform DrugBAN, further confirming the effectiveness of the FusionDTI framework in addressing cross-domain prediction challenges. 433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

These findings highlight not only the substantial improvements of FusionDTI over existing approaches but also its effectiveness in capturing finegrained information on DTI. The key to this success lies in FusionDTI's token-level fusion module, which enables the model to consider fine-grained interactions for each drug-target pair. This finegrained interaction information aligns closely with biomedical pathways, where binding events often depend on the specific atoms or substructures involved in interactions with residues. Therefore, the model's ability to capture such fine-grained interactions significantly enhances its predictive performance for DTI.



Figure 3: Performance comparison of two fusion strategies: BAN and CAN on the BindingDB.

CAN	AUC	AUPRC	Accuracy
×	0.954	0.963	0.894
√	0.989	0.990	0.961

Table 3: Ablation study of the CAN module on theBindingDB dataset.

4.3 Comparison of the BAN and CAN

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

There are two fusion strategies available: BAN and CAN, thus determining which one works better is a key step for establishing FusionDTI's prediction effectiveness. We perform a fair comparison involving the same encoders, classifier and dataset. As shown in Figure 3, we compare BAN and CAN by employing two linear layers to adjust the feature dimensions of the drug and target representations. With the feature dimension increasing, the performance of FusionDTI-CAN continues to rise, while that of FusionDTI-BAN reaches a plateau. When the feature dimension is 512, both of the variants attain their peak positions with an AUC of 0.989 and 0.967, respectively. These results indicate that the CAN module seems to be better suited to the DTI prediction tasks and in capturing fine-grained interaction information. In contrast, BAN may not be able to fully capture fine-grained binding information between proteins and drugs, such as the specific interactions between the drug atoms and residues. Therefore, these findings suggest that the CAN strategy is more effective and adaptable to the complexities involved in DTI prediction, providing superior performance, especially as the feature dimension scales.

4.4 Ablation Study

The fine-grained interaction of drug and target representations is critical in DTI as it directly impacts



Figure 4: Performance evaluation of fusion scales on the BindingDB dataset.

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

the model's ability to infer potential binding sites. For FusionDTI, this interaction is facilitated by the CAN module, which markedly enhances the predictive accuracy by capturing the fine-grained interaction information between the drugs and targets. Table 3 demonstrates the impact of the CAN module on the prediction performance. When the fusion module is omitted, the model achieves an AUC of 0.954 and an accuracy of 0.894. Conversely, using the CAN module, there is a significant improvement, with the AUC increasing to 0.989 and the accuracy reaching 0.961. This highlights the effectiveness of the CAN module in improving the inference ability of FusionDTI. In Appendix A.7 and A.8, we further compare time-consuming and time complexity with baselines.

4.5 Analysis of Fusion Scales

In assessing fusion representations, it is critical to determine whether more fine-grained modelling enhances the predictive performance. Thus, we define a grouping function with the parameter g (Group size) for averaging tokens within each group before the CAN fusion module. The parameter g, representing the number of tokens per group, controls the granularity of the attention mechanism. Specifically, when g is set to 1, the fusion operates at the token level, where each token is considered independently. In contrast, when \mathbf{g} is set to 512, the fusion occurs at a global level, considering the entire embedding as a single unit. We have the flexibility to control the fusion scale for the drug and protein representations, but the token length must be divisible by the group size. As shown in Figure 4, as the number of tokens per group increases from 1 to 512 (Maximum Token Length), the performance of the FusionDTI model declines

Drug-Target Interactions
 EZL - 6QL2: 1. sulfonamide oxygen - Leu198, Thr199 and Trp209; 2. amino group - His94, His96, His119 and Thr199; 3. benzothiazole ring - Leu198, Thr200, Tyr131, Pro20 and Gln92; 4. ethoxy group - Gln135;
 9YA - 5W8L: 1. amino group of sulfonamide - Asp140, Glu191; 2. sulfonamide oxygen - Asp140, Ile141 and Val139; 3. carboxylic acid oxygens - Arg168, His192, Asp194 ar Thr247; 4. biphenyl rings - Arg105, Asn137 and Pro138; 5. hydrophobic contact - Ala237, Tyr238 and Leu322;
 EJ4 - 4N6H: 1. basic nitrogen of ligand - Asp128; 2. hydrophobic pocket - Tyr308, Ile304 and Tyr129; 3. water molecules - Tyr129, Met132, Trp274, Tyr308

Table 4: FusionDTI predictions: **Bold** represents new predictions versus DrugBAN.

accordingly. This also aligns with the biomedical rules governing drug-protein interactions, where the principal factor influencing the binding is the interplay between the key atoms or substructures in the drug and primary residues in the protein. Furthermore, the CAN module outperforms BAN consistently at various scale settings, indicating that CAN better accesses the information between the drug and target. Consequently, this supports that the more detailed the interaction information obtained between the drugs and targets by the fusion module, the more beneficial it is for the enhancement of the model's prediction performance.

4.6 Case Study

and Lys214;

520

521

522

524

527

529

530

531

532

533

534

535

537

539

540

546

547

550

A further strength of FusionDTI to enable explainability, which is critical for drug design efforts, is the visualisation of each token's contribution to the final prediction through cross-attention maps. To compare with the DrugBAN model, we examine three identical pairs of DTI from the Protein Data Bank (PDB) (Berman et al., 2007): (EZL -6QL2 (Kazokaitė et al., 2019), 9YA - 5W8L (Rai et al., 2017) and EJ4 - 4N6H (Fenalti et al., 2014)), which are excluded from the training data. As shown in Table 4, our proposed model predicts more binding sites existing in the PDB (Berman et al., 2007) (in bold) by ranking the binding sites shown in the attention map. For instance, to predict the interaction of the drug EZL with the target 6QL2, our proposed model using BertViz (Vig, 2019) highlights potential binding sites as illus-



Figure 5: EZL - 6QL2: Fine-grained interactions via attention visualization.

551

552

553

554

555

556

557

558

559

561

562

563

564

565

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

trated in Figure 5. Specifically, our CAN module is effective in capturing fine-grained binding information at the token level, as we have successfully predicted the novel binding between Gln92 and the benzothiazole ring (Di Fiore et al., 2008). In particular, we address the lack of structural information on protein sequences by employing the SA vocabulary, which matches each residue to a corresponding 3D feature via Foldseek (Van Kempen et al., 2024). This study highlights the effectiveness of FusionDTI in enhancing performance on the DTI task, thereby supporting more targeted and efficient drug development efforts. In Appendix A.9, we further investigate ten DTI pairs in non-small cell lung cancer (NSCLC) from PDB (Waliany et al., 2025), highlighting predicted binding residues.

5 Conclusions

With the rapid increase of new diseases and the urgent need for innovative drugs, it is critical to capture fine-grained interactions, since the binding of specific drug atoms to the main amino acids is key to the DTI task. Despite some achievements, fine-grained interaction information is not effectively captured. To address this challenge, we introduce FusionDTI uses token-level fusion to effectively obtain fine-grained interaction information. Through experiments on three well-known datasets, we demonstrate that our proposed FusionDTI model outperforms eight state-of-the-art baselines, particularly in the more realistic crossdomain scenario. Additionally, we show that the attention weights of the token-level fusion module can highlight potential binding sites, providing a certain level of explainability.

585 Limitations

Even if our proposed model identifies potentially 586 useful DTI, these predictions need to be validated 587 by wet experiments, a time-consuming and expen-588 sive process. We have shown that FusionDTI is effective and efficient in screening for possible DTI in large-scale data as well as in locating potential binding sites in the process of drug design. How-592 ever, it is not directly applicable to human medical therapy and other biomedical interactions because it lacks clinical validation and regulatory approval for medical use. 596

References

599

601

602

606

610

611

612

613

614

615

617

618

621

- Walid Ahmad, Elana Simon, Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. 2022.
 Chemberta-2: Towards chemical foundation models. <u>arXiv preprint arXiv:2209.01712</u>.
 - Heba Askr, Enas Elgeldawi, Heba Aboul Ella, Yaseen AMM Elshaier, Mamdouh M Gomaa, and Aboul Ella Hassanien. 2023. Deep learning in drug discovery: an integrative review and future challenges. <u>Artificial Intelligence Review</u>, 56(7):5975–6037.
 - Peizhen Bai, Filip Miljković, Bino John, and Haiping Lu. 2023. Interpretable bilinear attention network with domain adaptation improves drugtarget prediction. <u>Nature Machine Intelligence</u>, 5(2):126–136.
 - Helen Berman, Kim Henrick, Haruki Nakamura, and John L Markley. 2007. The worldwide protein data bank (wwpdb): ensuring a single, uniform archive of pdb data. <u>Nucleic acids research</u>, 35(suppl_1):D301–D303.
- Ashraf Brik and Chi-Huey Wong. 2003. Hiv-1 protease: mechanism and drug discovery. Organic & biomolecular chemistry, 1(1):5–14.
- Dong-Sheng Cao, Qing-Song Xu, and Yi-Zeng Liang. 2013. propy: a tool to generate various modes of chou's pseaac. <u>Bioinformatics</u>, 29(7):960–962.
- Ashish Chandwani and Jonathan Shuter. 2008.
 Lopinavir/ritonavir in the treatment of hiv-1 infection: a review. <u>Therapeutics and clinical risk</u>
 <u>management</u>, 4(5):1023–1033.

Lifan Chen, Xiaoqin Tan, Dingyan Wang, Feisheng Zhong, Xiaohong Liu, Tianbiao Yang, Xiaomin Luo, Kaixian Chen, Hualiang Jiang, and Mingyue Zheng. 2020. Transformercpi: improving compound–protein interaction prediction by sequence-based deep learning with selfattention mechanism and label reversal experiments. <u>Bioinformatics</u>, 36(16):4406–4414. 629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

- Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. <u>Machine learning</u>, 20:273–297.
- Mindy I Davis, Jeremy P Hunt, Sanna Herrgard, Pietro Ciceri, Lisa M Wodicka, Gabriel Pallares, Michael Hocker, Daniel K Treiber, and Patrick P Zarrinkar. 2011. Comprehensive analysis of kinase inhibitor selectivity. <u>Nature biotechnology</u>, 29(11):1046–1051.
- Anna Di Fiore, Carlo Pedone, Jochen Antel, Harald Waldeck, Andreas Witte, Michael Wurl, Andrea Scozzafava, Claudiu T Supuran, and Giuseppina De Simone. 2008. Carbonic anhydrase inhibitors: the x-ray crystal structure of ethoxzolamide complexed to human isoform ii reveals the importance of thr200 and gln92 for obtaining tight-binding inhibitors. <u>Bioorganic &</u> medicinal chemistry letters, 18(8):2669–2674.
- Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Wang Yu, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, Debsindhu Bhowmik, and Burkhard Rost. 2021. Prottrans: Towards cracking the language of lifes code through selfsupervised deep learning and high performance computing. <u>IEEE Transactions on Pattern</u> Analysis and Machine Intelligence, pages 1–1.
- Gustavo Fenalti, Patrick M Giguere, Vsevolod Katritch, Xi-Ping Huang, Aaron A Thompson, Vadim Cherezov, Bryan L Roth, and Raymond C Stevens. 2014. Molecular control of δ -opioid receptor signalling. <u>Nature</u>, 506(7487):191–196.
- Bowen Gao, Bo Qiang, Haichuan Tan, Yinjun Jia, Minsi Ren, Minsi Lu, Jingjing Liu, Wei-Ying Ma, and Yanyan Lan. 2023. Drugclip: Contrastive protein-molecule representation learning for virtual screening. <u>Advances in Neural</u> <u>Information Processing Systems</u>, 36:44595– 44614.

- 723 724 725 727 728 729 730 732 733 734 735 736 737 738 739 740 Unsupervised sentence-pair 741 742 743 744 745 Improv-746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768
- Michael K Gilson, Tiqing Liu, Michael Baitaluk, 676 George Nicola, Linda Hwang, and Jenny Chong. 677 2016. Bindingdb in 2015: a public database for medicinal chemistry, computational chem-679 istry and systems pharmacology. Nucleic acids research, 44(D1):D1045-D1053.
 - Yichen Gong, Heng Luo, and Jian Zhang. 2018. Natural language inference over interaction space. International Conference on Learning Representations.
 - Mercedes Herrera-Juárez, Cristina Serrano-Gómez, Helena Bote-de Cabo, and Luis Paz-Ares. 2023. Targeted therapy for lung cancer: Beyond egfr and alk. Cancer, 129(12):1803-1820.
 - Tin Kam Ho. 1995. Random decision forests. In Proceedings of 3rd international conference on document analysis and recognition, volume 1, pages 278-282. IEEE.
 - Kexin Huang, Cao Xiao, Lucas M Glass, and Jimeng Sun. 2021. Moltrans: molecular interaction transformer for drug-target interaction prediction. Bioinformatics, 37(6):830-836.

701

706

707

711

713

714

715

716

717

718

719

721

- Justina Kazokaitė, Visvaldas Kairys, Joana Smirnovienė, Alexey Smirnov, Elena Manakova, Martti Tolvanen, Seppo Parkkila, and Daumantas Matulis. 2019. Engineered carbonic anhydrase vi-mimic enzyme switched the structure and affinities of inhibitors. Scientific reports, 9(1):12710.
- Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Bilinear attention networks. Zhang. 2018. Advances in neural information processing systems, 31.
- Mario Krenn, Qianxiang Ai, Senja Barthel, Nessa Carson, Angelo Frei, Nathan C Frey, Pascal Friederich, Théophile Gaudin, Alberto Alexander Gayle, Kevin Maik Jablonka, and 1 others. 2022. Selfies and the future of molecular string representations. Patterns, 3(10).
- Ingoo Lee, Jongsoo Keum, and Hojung Nam. 2019. Deepconv-dti: Prediction of drug-target interactions via deep learning with convolution on protein sequences. PLoS computational biology, 15(6):e1007129.
 - Peizhao Li, Jiuxiang Gu, Jason Kuen, Vlad I Morariu, Handong Zhao, Rajiv Jain, Varun Man-

junatha, and Hongfu Liu. 2021. Selfdoc: Selfsupervised document representation learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5652-5660.

- Shuqi Li, Shufang Xie, Hongda Sun, Yuhan Chen, Tao Qin, Tianjun Ke, and Rui Yan. 2025. Min: Multi-channel interaction network for drug-target interaction with protein distillation. IEEE Transactions on Computational Biology and Bioinformatics.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, and 1 others. 2023. Evolutionary-scale prediction of atomic-level protein structure with a language model. Science, 379(6637):1123-1130.
- Fangyu Liu, Yunlong Jiao, Jordan Massiah, Emine Yilmaz, and Serhii Havrylov. 2021. Trans-encoder: modelling through self-and mutual-distillations. In International Conference on Learning Representations.
- Hui Liu, Jianjiang Sun, Jihong Guan, Jie Zheng, and Shuigeng Zhou. 2015. ing compound-protein interaction prediction by building up highly credible negative samples. Bioinformatics, 31(12):i221-i229.
- Michael M Mysinger, Michael Carchia, John J Irwin, and Brian K Shoichet. 2012. Directory of useful decoys, enhanced (dud-e): better ligands and decoys for better benchmarking. Journal of medicinal chemistry, 55(14):6582-6594.
- Ha Cong Nga, Phuc Pham, and Truong Son Hy. 2025. Lantern: Leveraging large language models and transformers for enhanced molecular interactions. bioRxiv, pages 2025-02.
- Thin Nguyen, Hang Le, Thomas P Quinn, Tri Nguyen, Thuc Duy Le, and Svetha Venkatesh. Graphdta: predicting drug-target 2021. binding affinity with graph neural networks. Bioinformatics, 37(8):1140-1147.
- Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. 2023. BioT5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. In Proceedings

769	of the 2023 Conference on Empirical Methods	2024. Saprothub: Making protein modeling ac-	815
770	in Natural Language Processing, pages 1102-	cessible to all biologists. bioRxiv, pages 2024–	816
771	1123, Singapore. Association for Computational	05.	817
772	Linguistics.	Von Sun, Von Vili Corcon VI aung and Dingzhoo	0.14
773	Lihong Peng, Xin Liu, Long Yang, Longlong Liu	Hu 2024 ingan dti: prediction of drug target	010
774	Zongzheng Bai Min Chen Xu Lu, and Libo	interaction with interpretable posted graph pay	000
775	Nie 2024 Bindti: A bi-directional intention	rel network and pretrained melacula medale	021
776	network for drug-target interaction identification	Picinformatics 40(2):http://25	82
777	based on attention mechanisms. IEEE Journal	$\underline{\text{Bioinformatics}}, 40(5).000155.$	824
778	of Biomedical and Health Informatics	Emma Svensson Pieter-Ian Hoedt Sepp Hochre-	823
110	<u>or Diomedical and Ticatal Informatics</u> .	iter, and Günter Klambauer, 2024. Hyper-	824
779	Ganesha Rai, Kyle R Brimacombe, Bryan T	pcm: Robust task-conditioned modeling of	82!
780	Mott, Daniel J Urban, Xin Hu, Shyh-Ming	drug-target interactions. Journal of Chemical	826
781	Yang, Tobie D Lee, Dorian M Cheff, Jennifer	Information and Modeling, 64(7):2539–2553.	827
782	Kouznetsova, Gloria A Benavides, and 1 others.		
783	2017. Discovery and optimization of potent, cell-	Jing Tang, Agnieszka Szwajda, Sushil Shakyawar,	828
784	active pyrazole-based inhibitors of lactate dehy-	Tao Xu, Petteri Hintsanen, Krister Wennerberg,	829
785	drogenase (ldh). Journal of medicinal chemistry,	and Tero Aittokallio. 2014. Making sense	830
786	60(22):9184–9204.	of large-scale kinase inhibitor bioactivity data	831
		sets: a comparative and integrative analysis.	832
787	David Rogers and Mathew Hahn. 2010. Extended-	Journal of chemical information and modeling,	833
788	connectivity fingerprints. Journal of chemical	54(3):735–743.	834
789	information and modeling, 50(5):742–754.		
790	Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie	Michel Van Kempen, Stephanie S Kim, Char-	835
791	Ying Wei Wenhing Huang and Junzhou Huang	lotte Tumescheit, Milot Mirdita, Jeongjae Lee,	836
792	2020 Self-supervised graph transformer on	Cameron LM Gilchrist, Johannes Söding, and	837
793	large-scale molecular data. Advances in neural	Martin Steinegger. 2024. Fast and accurate	838
794	information processing systems, 33:12559–	protein structure search with foldseek. <u>Nature</u>	839
795	12571.	<u>Biotechnology</u> , $42(2):243-246$.	84(
		Mihaly Varadi Stephen Anyango Mandar Desh	0./-
796	Jerret Ross, Brian Belgodere, Vijil Chenthamarak-	nande Sreenath Nair Cindy Natassia Gal-	04
797	shan, Inkit Padhi, Youssef Mroueh, and Payel	abina Vordanova David Vuan Oana Stroe	8/12
798	Das. 2022. Large-scale chemical language repre-	Gemma Wood Agata Laydon and 1 others	844
799	sentations capture molecular structure and prop-	2022 Alphafold protein structure database:	84
800	erties. <u>Nature Machine Intelligence</u> , 4(12):1256–	massively expanding the structural coverage	846
801	1264.	of protein-sequence space with high-accuracy	847
202	Monica Schenone, Vlado Dančík, Bridget K Wag-	models. Nucleic acids research, 50(D1):D439–	848
803	ner and Paul A Clemons 2013 Target identi-	D444.	849
804	fication and mechanism of action in chemical		
805	biology and drug discovery Nature chemical	Ashish Vaswani, Noam Shazeer, Niki Parmar,	850
806	hiology 9(4):232–240	Jakob Uszkoreit, Llion Jones, Aidan N Gomez,	851
000	<u>01010gj</u> , y(1).202 210.	Łukasz Kaiser, and Illia Polosukhin. 2017. At-	852
807	Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan,	tention is all you need. Advances in neural	853
808	Xibin Zhou, and Fajie Yuan. 2023. Saprot: pro-	information processing systems, 30.	854
809	tein language modeling with structure-aware		
810	vocabulary. Advances in neural information	Jesse Vig. 2019. A multiscale visualization	855
811	processing systems, pages 2023-10.	of attention in the transformer model. In	856
		Proceedings of the 57th Annual Meeting of	857
812	Jin Su, Zhikai Li, Chenchen Han, Yuyang Zhou,	the Association for Computational Linguistics:	858
813	Junjie Shan, Xibin Zhou, Dacheng Ma, The	System Demonstrations, pages 37–42, Florence,	859
814	OPMC, Sergey Ovchinnikov, and Fajie Yuan.	Italy. Association for Computational Linguistics.	860

- Sarah Waliany, Jessica J Lin, and Justin F Gainor.
 2025. Evolution of first versus next-line targeted
 therapies for metastatic non-small cell lung cancer.
 <u>Trends in Cancer</u>.
- 865David Weininger. 1988. Smiles, a chemical lan-
guage and information system. 1. introduction
to methodology and encoding rules. Journal
of chemical information and computer sciences,
28(1):31–36.

871

872

874

875

876

877

878

879

892

893

894

899

900

- David Weininger, Arthur Weininger, and Joseph L Weininger. 1989. Smiles. 2. algorithm for generation of unique smiles notation. Journal of chemical information and computer sciences, 29(2):97–101.
- David S Wishart, Craig Knox, An Chi Guo, Dean Cheng, Savita Shrivastava, Dan Tzur, Bijaya Gautam, and Murtaza Hassanali. 2008. Drugbank: a knowledgebase for drugs, drug actions and drug targets. <u>Nucleic acids research</u>, 36(suppl_1):D901–D906.
 - Yifan Wu, Min Gao, Min Zeng, Jie Zhang, and Min Li. 2022. Bridgedpi: a novel graph neural network for predicting drug–protein interactions. <u>Bioinformatics</u>, 38(9):2571–2578.
 - Minghao Xu, Xinyu Yuan, Santiago Miret, and Jian Tang. 2023. Protst: Multi-modality learning of protein sequences and biomedical texts. In <u>International Conference on Machine Learning</u>, pages 38749–38767. PMLR.
 - Ziduo Yang, Weihe Zhong, Lu Zhao, and Calvin Yu-Chian Chen. 2021. Ml-dti: mutual learning mechanism for interpretable drug-target interaction prediction. <u>The Journal of Physical</u> <u>Chemistry Letters</u>, 12(17):4247–4261.
 - Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. 2021. Do transformers really perform badly for graph representation? <u>Advances in neural information</u> processing systems, 34:28877–28888.
- 901Atakan Yüksel, Erva Ulusoy, Atabey Ünlü, and
Tunca Doğan. 2023. Selformer: molecular rep-
resentation learning via selfies language models.903Machine Learning: Science and Technology,
4(2):025035.

Dataset	Drugs	Proteins	Interactions
BindingDB	14,643	2,623	49,199
BioSNAP	4,510	2,181	27,464
Human	2,726	2,001	6,728
DAVIS	68	442	30,056
KIBA	2,068	229	118,254
DUD-E	1,200,966	102	1,434,019

Table 5:	Dataset	Statistics.
----------	---------	-------------

Xiaoting Zeng, Weilin Chen, and Baiying Lei. 906 2024. Cat-dti: cross-attention and transformer 907 network with domain adaptation for drug-target 908 interaction prediction. BMC bioinformatics, 909 25(1):141. 910 Hongzhi Zhang, Xiuwen Gong, Shirui Pan, Jia 911 Wu, Bo Du, and Wenbin Hu. 2024. A cross-912 field fusion strategy for drug-target interaction 913 prediction. arXiv preprint arXiv:2405.14545. 914 Shuangjia Zheng, Yongjian Li, Sheng Chen, Jun 915 Xu, and Yuedong Yang. 2020. Predicting drug-916 protein interaction using quasi-visual question 917 answering system. Nature Machine Intelligence, 918 2(2):134-140. 919 Marinka Zitnik, Rok Sosic, and Jure Leskovec. 920 2018. Biosnap datasets: Stanford biomedical 921 network dataset collection. Note: http://snap. 922 stanford. edu/biodata Cited by, 5(1). 923 Appendix Α 924 A.1 Hyperparameter of FusionDTI 925 FusionDTI is implemented in Python 3.8 and the 926 PyTorch framework $(1.12.1)^2$. The computing de-927 vice we use is the NVIDIA GeForce RTX 3090. 928 929

In the "Experimental Setup and Results" section, we only present experiment results based on the BindingDB dataset, as the performance trends are identical to the BioSNAP dataset and the Human dataset. Table 6 shows the parameters of the FusionDTI model and Table 7 lists the notations used in this paper with descriptions.

930

931

932

933

934

935

936

937

938

939

A.2 Dataset Sources

All the data used in this paper are from public sources. The statistics of the experimental datasets are presented in Table 5.

²https://pytorch.org/

Module	Hyperparameter	Value
Mini-batch	Batch size	64 (options: 64, 128)
Drug Encoder	PLM	HUBioDataLab/SELFormer
Protein Encoder	PLM	westlake-repl/SaProt_650M_AF2
BAN	Heads of bilinear attention	3
	Bilinear embedding size	512 (options: 32, 64, 128, 256, 512, 768)
	Sum pooling window size	2
CAN	Attention heads	8
	Hidden dimension	512 (options: 32, 64, 128, 256, 512, 768)
	Integration strategies	Mean pooling (options: Mean pooling, CLS)
	Group size	1 (options: from 1 to 512)
MLP	Hidden layer sizes	(1024, 512, 256)
	Activation	Relu (options: Tanh, Relu)
	Solver	AdamW
		(options: AdamW, Adam, RMSprop, Adadelta, LBFGS)
	Learning rate scheduler	CosineAnnealingLR
		(options: CosineAnnealingLR, StepLR, ExponentialLR)
	Initial learning rate	1e-4 (options: from 1e-3 to 1e-6)
	Maximum epoch	200

Table 6: Configuration Parameters

Notations	Description
D	Drug feature
Р	Target feature
$\mathbf{q} \in \mathbb{R}^{K}$	weight vector for bilinear transformation
$Att \in \mathbb{R}^{\rho \times \phi}$	Bilinear attention maps in BAN
$\mathbf{U} \in \mathbb{R}^{N imes K}$	Transformation matrix for drug features
$\mathbf{V} \in \mathbb{R}^{M imes K}$	Transformation matrix for target features
g	The number of tokens per group
$\mathbf{D}^* \in \mathbb{R}^{m imes h}$	Fused drug representations in token-level interaction
$\mathbf{P}^* \in \mathbb{R}^{n imes h}$	Fused target representations in token-level interaction
$\mathbf{Q}_d, \mathbf{K}_d, \mathbf{V}_d \in \mathbb{R}^{m imes h}$	Queries, keys, and values for the drug in token-level interaction
$\mathbf{Q}_p, \mathbf{K}_p, \mathbf{V}_p \in \mathbb{R}^{n imes h}$	Queries, keys, and values for target in token-level interaction
$\mathbf{W}_{q}^{d}, \mathbf{W}_{k}^{d}, \mathbf{W}_{v}^{d} \in \mathbb{R}^{H \times h}$	Projection matrices for drug queries, keys, and values
$\mathbf{W}_{q}^{\dot{p}},\mathbf{W}_{k}^{p},\mathbf{W}_{v}^{p}\in\mathbb{R}^{h imes h}$	Projection matrices for target queries, keys, and values
\mathbf{F}	drug-target joint representation
$p \in [0, 1]$	output interaction probability
H	Number of attention heads in token-level interaction
m, n	Sequence lengths for drug and protein respectively
h	Hidden dimension in token-level interaction

Table 7: Notations and Descriptions

 The BindingDB (Gilson et al., 2016) dataset
 is a web-accessible database of experimentally validated binding affinities, focusing
 primarily on the interactions of small druglike molecules and proteins. The BindingDB
 source is found at https://www.bindingdb.

org/bind/index.jsp.

946

2. The BioSNAP (Zitnik et al., 2018) dataset is created from the DrugBank database (Wishart et al., 2008). It is a balanced dataset with validated positive interactions and an equal number of negative samples randomly ob-

1001

tained from unseen pairs. The BioSNAP source is found at https://github.com/kexinhuang12345/MolTrans.

952

953

955

961

962

963

964

965

966

967

969

970

971

972

973

974

975

977

980

981

982

985

987

988

991

993

995

996

997

999

- 3. The Human (Liu et al., 2015; Chen et al., 2020) dataset includes highly credible negative samples. The balanced version of the Human dataset contains the same number of positive and negative samples. The Human source is found at https://github.com/ lifanchen-simm/transformerCPI.
- 4. The DAVIS (Davis et al., 2011) dataset provides continuous binding affinity measurements (K_d values) between kinase inhibitors and proteins. It is commonly used for regression-based drug-target interaction (DTI) prediction tasks. The DAVIS source is available at https://tdcommons. ai/multi_pred_tasks/dti/.
 - 5. The KIBA (Tang et al., 2014) dataset integrates multiple bioactivity measures to provide a unified KIBA score for kinase-inhibitor pairs. It is widely adopted in benchmark studies for affinity prediction. The KIBA source is available at https://tdcommons. ai/multi_pred_tasks/dti/.
 - 6. The DUD-E (Mysinger et al., 2012) (Directory of Useful Decoys, Enhanced) dataset is a large-scale benchmark set for virtual screening, containing active compounds and challenging decoys for various protein targets. The DUD-E source is found at http://dude.docking.org/.

A.3 How to Obtain the Structure-aware (SA) Sequence of a Protein and the SELFIES of a Drug?

To obtain the SA sequence of a protein, the first step is to obtain Uniprot IDs from the UniProt website using information such as the amino acid sequences or protein names, and then save these IDs in a comma-delimited text file. Subsequently, we use the UniProt IDs to fetch the relevant 3D structure file (.cif) from AlphafoldDB (Varadi et al., 2022) using Foldseek. The SA vocabulary of the protein can then be generated from this 3D structure file.

For drugs, the SELFIES could be derived from SMILES strings. This conversion requires specific Python packages, and upon installation, the SELF-IES strings can be generated through appropriate scripts. Please refer to our submission file for detailed procedures, including the necessary code.

Notably, our submission of supplementary material contains step-by-step descriptions and code for generating the SA sequences and SELFIES.

A.4 Baselines

We compare the performance of FusionDTI with the following eight models on the DTI task.

Baselines on BindingDB, BioSNAP, and Human.

- Support Vector Machine (Cortes and Vapnik, 1995) on the concatenated fingerprint ECFP4 (Rogers and Hahn, 2010) (extended connectivity fingerprint, up to four bonds) and PSC (Cao et al., 2013) (pseudo-amino acid composition) features.
- 2. Random Forest (Ho, 1995) on the concatenated fingerprint ECFP4 and PSC features.
- 3. DeepConv-DTI (Lee et al., 2019) uses a fully connected neural network to encode the ECFP4 drug fingerprint and a CNN along with a global max-pooling layer to extract features from the protein sequences. Then the drug and protein features are concatenated and fed into a fully connected neural network for the final prediction.
- 4. GraphDTA (Nguyen et al., 2021) uses GNN for the encoding of drug molecular graphs, and a CNN is used for the encoding of the protein sequences. The derived vectors of the drug and protein representations are directly concatenated for interaction prediction.
- 5. MolTrans (Huang et al., 2021) uses a transformer architecture to encode the drugs and proteins. Then a CNN-based fusion module is adapted to capture DTI interactions.
- 6. DrugBAN (Bai et al., 2023) use a Graph Convolution Network and 1D CNN to encode the drug and protein sequences. Then a bilinear attention network (Kim et al., 2018) is adopted to learn pairwise interactions between the drug and protein. The resulting joint representation is decoded by a fully connected neural network.
- 7. BioT5 (Pei et al., 2023) is a cross-modeling model in biology with chemical knowledge and natural language associations.

- 1046 1047
- 1048 1049
- 1050
- 1052
- 1053
- 1054 1055
- 1056 1057
- 1058
- 1059
- 1060 1061
- 1062 1063
- 1065
- 1067
- 1068
- 1070
- 1071
- 1072 1073
- 1074
- 1075 1076
- 1077 1078
- 1079 1080
- 1082

- 1085 1086
- 1084
- 1087 1088

1092



Baselines on DAVIS and KIBA.

- 9. ML-DTI (Yang et al., 2021) combines molecular fingerprints with physicochemical descriptors and applies MLPs for regression.
- 10. DGraphDTA (Alphafold2) (Wu et al., 2022) integrates protein 3D structural data (from AlphaFold2) with drug graphs through a dualgraph encoding strategy.
- 11. iNGNN-DTI (Sun et al., 2024) introduces an interpretable graph neural network with attention-based gating mechanisms for drug-target regression tasks.
- 12. MIN (Li et al., 2025) uses a hierarchical multichannel network that combines structureaware and structure-agnostic representations with interpretable attention mechanisms.

Baselines on DUD-E.

- 13. DrugVQA (Zheng et al., 2020) formulates DTI prediction as a visual question answering task over molecular structures and protein sequences.
 - 14. DrugClip (Gao et al., 2023) adapts a contrastive pretraining framework, aligning drug molecules and protein embeddings using a CLIP-style architecture.
 - 15. HyperPCM (Svensson et al., 2024) utilises hyperbolic protein-compound matching for robust generalisation in few-shot virtual screening scenarios.
- 16. MIN (Li et al., 2025) introduces multiinstance networks to model DTI at the binding site level using hierarchical attention.

A.5 Ablation Study

In Table 8, we compare the performance of two aggregation strategies within the CAN module. The pooling strategy outperforms the CLS-based aggregation, achieving an AUC and AUPRC of 0.989 and 0.990, respectively. This comparison highlights the superior effectiveness of the pooling in aggregating contextual information. Thus, the integration of a CAN module, particularly employing a pooling aggregation strategy, is shown to be essential for making confident and accurate predictions.



Figure 6: Time comparison on the BindingDB, Human and BioSNAP datasets.

Aggregation	AUC	AUPRC	Accuracy
CLS	0.982	0.983	0.956
Pooling	0.989	0.990	0.961

Table 8: Comparison of aggregation strategies for FusionDTI-CAN on the BindingDB dataset.

1093

1116

Evaluation of PLMs Encoding A.6

The protein encoder and drug encoder are funda-1094 mental for the token-level fusion of representa-1095 tions, as these encoders are responsible for gen-1096 erating fine-grained representations to better explore interaction information. Our proposed model 1098 employs two PLMs encoding two biomedical en-1099 tities: the drug and protein, respectively. In 1100 terms of the protein encoders, Figure 7 com-1101 pares the the performance of the two protein en-1102 coders (SaProt (Su et al., 2023) and ESM-2 (Lin 1103 et al., 2023)) in combination with three differ-1104 ent drug encoders: ChemBERTa-2 (Ahmad et al., 1105 2022), SELFormer (Yüksel et al., 2023) and MoL-1106 Former (Ross et al., 2022). From the figure, we find 1107 that SaProt consistently outperforms ESM-2 when 1108 combined with all three drug encoders. As can be 1109 seen in Figure 8, SELFormer achieves the best per-1110 formance in encoding the drug sequences among 1111 the three advanced drug encoders. Notably, the top-1112 performing combination is SaProt and SELFormer, 1113 hence our proposed FusionDTI uses them as drug 1114 and protein encoders. 1115

A.7 **Efficiency Analysis**

Efficiency in computational models is crucial, par-1117 ticularly when handling large-scale and exten-1118 sive datasets in drug discovery. Our proposed 1119 model stores drug representations and target rep-1120 resentations in memory for later online training. 1121



Figure 7: Performance comparison of protein encoders on the BindingDB dataset.

1122 As evidenced by Figure 6, FusionDTI-CAN and FusionDTI-BAN with pre-encoded representations 1123 process the BindingDB dataset much faster than the 1124 non-pre-coded models, approximately 45 minutes 1125 and 220 minutes, respectively. This stark difference 1126 highlights the advantage of pre-encoded, which 1127 eliminates the need for real-time data processing 1128 and accelerates the overall throughput. While 1129 FusionDTI-BAN and DrugBAN have the same 1130 fusion module, the pre-encoded FusionDTI-BAN 1131 runs faster and predicts more accurately, as shown 1132 in Table 1. In addition, FusionDTI-BAN runs 1133 faster than FusionDTI-CAN, indicating that the 1134 BAN fusion module is more efficient. Ultimately, 1135 FusionDTI-BAN with pre-encoded data stands out 1136 as a highly efficient approach, offering substantial 1137 benefits in scenarios where exists large-scale data. 1138

A.8 Time Complexity Analysis

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

Fusion module	Complexity (O)	Parameters
BAN	$O(\rho \cdot \phi \cdot K)$	790k
CAN	$O(m \cdot n \cdot h)$	1572k

Table 9: Time complexity and parameters comparison of BAN and CAN.

The feature dimensions of the representations generated by different PLM encoders are fixed, but the size of the feature dimensions may not be the same. Therefore, in order to fuse protein and drug representations, we use two linear layers to keep the representations' feature dimension equal to the token length (512).

The time complexity of BAN depends on the computation of bilinear interaction maps. The



Figure 8: Performance comparison of drug encoders on the BindingDB dataset.

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

bilinear attention involves a Hadamard product and further matrix operations as given in Equation (2). The computation of $U^T P$ and $V^T D$ requires $O(N \cdot \rho \cdot K)$ and $O(M \cdot \phi \cdot K)$ operations, respectively. Here, K denotes the dimensionality of the transformation, which is the rank of the feature space to which the protein and drug features are projected. When the token length is equal to the feature dimension and the dimensions of transformation are two times either, the overall time complexity is $O(\rho \cdot \phi \cdot K)$.

For the token-level interaction in the DTI task, the time complexity is also markedly influenced by the attention mechanisms. It also satisfies the condition that the token length is equal to the feature dimension of the drug and protein. With multi-head attention heads (H = 8), the complexity for computing the queries, keys, and values in the Equation (6) and (7), as well as the softmax attention weights, is given by $O(H \cdot n \cdot m \cdot h)$, where mandn represents the token lengths for the drug and protein, respectively, and h is the hidden dimension. Since each head contributes its own set of computations and the attention mechanism operates over all tokens, the $m \cdot n$ term (stemming from the softmax operation across the token length) becomes significant. This leads to a total time complexity of $O(m \cdot n \cdot h)$ per batch for the attention mechanism.

From the above analysis of the time complexity 1177 of the two fusion strategies, the time complexity of 1178 CAN is lower than BAN in the case of the same 1179 input protein and drug features. BAN is markedly 1180 affected by the transformation dimension K. When 1181 the K is larger than the token and feature dimen-1182 sion, the time complexity of BAN is higher than 1183 CAN. However, we observe that the number of pa-1184 rameters in BAN is smaller than that of CAN via 1185

Drug-Target (Ligand - PDB ID)	Predicted Binding Residues
VGH - 2YFX	Glu113, Val46, Gly117, Met115, Asp186, Arg125, Lys225, Gln50, Ala190, Pro319
C6F - 6JQR	Tyr126, Asp209, Ala72, Glu208, Glu197, Leu219, Pro163, Gln97, Val225, His151
5P8 - 4CLI	GLu113, Leu172, Gly118, Ala64, Asp186, Ala150, Ile99, Pro290, Ala312, Glu316
0WM - 4G5J	His296, Pro102, Pro156, Met295, Asn116, Ser92, Thr217, Lys237, His143, Trp188
YY3 - 6LUD	Phe102, Leu151, Met1100, Lys52, Glu111, Ile22, Pro60, Ala129, Val141, Gly42
AQ4 - 1M17	Leu155, Leu99, Met104, Phe106, Thr165, Asp111, Lys171, Trp209, Ala61, Asp280
YMX - 5FTO	Asn162, Gly110, Phe35, Glu118, Val38, His155, ALa197, Met46, Leu112, Asp280
1C9 - 4I23	Ala50, Leu95, Met100, Pro101, Glu69, Thr247, Tyr120, His177, Pro221, Val49
VGH - 2XP2	Leu172, Gly185, Ala116, Lys66, Asp119, Pro58, Met82, Pro131, Ala167, Val27
EMH - 3AOX	Glu143, Leu55, Gly56, Val113, Met132, Glu91, Leu157, Val44, Ala59, Ile166

Table 10: Predicted binding sites for DTI in NSCLC. **Bold** residues are supported by the PDB database, while others remain unverified.

the Pytroch package, as shown in Table 9.

A.9 Case Study

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1206

1207

1208

1209

1210

1211

The top three predictions (PDB ID: 6QL2 (Kazokaitė et al., 2019), 5W8L (Rai et al., 2017) and 4N6H (Fenalti et al., 2014)) of the cocrystalized ligands are derived from Protein Data Bank (PDB) (Berman et al., 2007). Following the setup of the DrugBAN case study, we only choose X-ray structures with a resolution greater than 2.5 Å corresponding to human proteins. In addition, the co-crystalized ligands are required to have pIC₅₀ \leq 100 nM and are not part of the training dataset.

To further DTI in non-small cell lung cancer (NSCLC), we identify ten additional drug-protein pairs from PDB. The selected targets—Epidermal Growth Factor Receptor (EGFR), Anaplastic Lymphoma Kinase (ALK), and ROS1—are wellestablished oncogenic drivers in NSCLC (Waliany et al., 2025). The corresponding inhibitors, including Erlotinib, Gefitinib, Osimertinib, Crizotinib, and Lorlatinib, exhibit high binding affinities (Herrera-Juárez et al., 2023). Table 10 presents the predicted binding residues for these interactions, with bolded residues supported by experimental PDB data, while others remain unverified.

A.10 Performance Comparison

1212Tables 11 and 12 provide a detailed performance
evaluation of FusionDTI and baseline models1213evaluation of FusionDTI and baseline models1214across both in-domain and cross-domain settings.1215To ensure a comprehensive assessment, we report1216multiple evaluation metrics, including AUROC and1217AUPRC as primary indicators, alongside F1-score,

Sensitivity, Specificity, and Matthews Correlation Coefficient (MCC). These additional metrics offer deeper insights into model performance across different classification aspects.

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

In addition, Tables 13, 14, and 15 present results on three benchmark datasets: DAVIS (Davis et al., 2011), KIBA (Tang et al., 2014), and DUD-E (Mysinger et al., 2012). Each table compares FusionDTI with strong task-specific baselines under standard evaluation metrics for their respective tasks, further demonstrating the robustness and adaptability of our model.

Model	AUC	AUPR	Accuracy	F1	Sensitivity	Specificity	MCC	
BindingDB								
SVM	$0.939 {\pm} 0.001$	$0.928 {\pm} 0.002$	0.825 ± 0.004	0.821 ± 0.004	$0.810{\pm}0.010$	$0.840 {\pm} 0.007$	$0.700{\pm}0.012$	
RF	$0.942{\pm}0.011$	$0.921 {\pm} 0.016$	$0.880{\pm}0.012$	$0.875 {\pm} 0.012$	$0.870 {\pm} 0.015$	$0.890 {\pm} 0.010$	$0.815 {\pm} 0.009$	
DeepConv-DTI	$0.945 {\pm} 0.002$	$0.925 {\pm} 0.005$	$0.882 {\pm} 0.007$	$0.878{\pm}0.008$	$0.870 {\pm} 0.011$	$0.885 {\pm} 0.010$	$0.818{\pm}0.013$	
GraphDTA	$0.951{\pm}0.002$	$0.934 {\pm} 0.002$	$0.888 {\pm} 0.005$	$0.884{\pm}0.005$	$0.880{\pm}0.006$	$0.890 {\pm} 0.004$	$0.825 {\pm} 0.008$	
MolTrans	$0.952{\pm}0.002$	$0.936{\pm}0.001$	$0.887 {\pm} 0.006$	$0.882{\pm}0.006$	$0.875 {\pm} 0.009$	$0.890 {\pm} 0.007$	$0.820{\pm}0.010$	
DrugBAN	$0.960 {\pm} 0.001$	$0.948 {\pm} 0.002$	$0.906 {\pm} 0.004$	$0.901 {\pm} 0.004$	$0.900 {\pm} 0.008$	$0.908 {\pm} 0.004$	$0.872 {\pm} 0.005$	
SiamDTI	0.961 ± 0.002	0.945 ± 0.002	$0.890 {\pm} 0.006$	$0.886 {\pm} 0.006$	$0.880 {\pm} 0.007$	$0.895 {\pm} 0.005$	0.830 ± 0.006	
BioT5	0.963 ± 0.001	$0.952 {\pm} 0.001$	0.907 ± 0.003	0.905 ± 0.003	0.900 ± 0.004	0.910 ± 0.003	0.850 ± 0.005	
FusionDTI-BAN	0.975 ± 0.002	$0.976 {\pm} 0.002$	$0.933 {\pm} 0.003$	$0.934{\pm}0.002$	$0.932 {\pm} 0.004$	0.935 ± 0.003	0.900 ± 0.003	
FusionDTI-CAN	$0.989{\pm}0.002$	$0.990{\pm}0.002$	$0.961{\pm}0.002$	$0.963{\pm}0.012$	$0.954{\pm}0.003$	$0.955{\pm}0.012$	$0.925{\pm}0.023$	
			BioSN	AP				
SVM	$0.862{\pm}0.007$	$0.864 {\pm} 0.004$	$0.777 {\pm} 0.011$	$0.773 {\pm} 0.011$	$0.760{\pm}0.015$	$0.780 {\pm} 0.008$	$0.690{\pm}0.013$	
RF	$0.860 {\pm} 0.005$	$0.886{\pm}0.005$	$0.804 {\pm} 0.005$	$0.800{\pm}0.005$	$0.795{\pm}0.008$	$0.810 {\pm} 0.007$	$0.715 {\pm} 0.006$	
DeepConv-DTI	$0.886{\pm}0.006$	$0.890{\pm}0.006$	$0.805 {\pm} 0.009$	$0.801 {\pm} 0.009$	$0.800{\pm}0.013$	$0.810 {\pm} 0.010$	$0.718{\pm}0.012$	
GraphDTA	$0.887 {\pm} 0.008$	$0.890{\pm}0.007$	$0.800 {\pm} 0.007$	$0.796 {\pm} 0.007$	$0.790{\pm}0.010$	$0.810 {\pm} 0.009$	$0.712{\pm}0.009$	
MolTrans	$0.895 {\pm} 0.004$	$0.897 {\pm} 0.005$	$0.825 {\pm} 0.010$	$0.820 {\pm} 0.010$	$0.815 {\pm} 0.013$	$0.830 {\pm} 0.012$	$0.730 {\pm} 0.011$	
DrugBAN	$0.903 {\pm} 0.005$	$0.902 {\pm} 0.004$	$0.834{\pm}0.008$	$0.830 {\pm} 0.009$	$0.820 {\pm} 0.021$	$0.847 {\pm} 0.010$	$0.719 {\pm} 0.007$	
SiamDTI	0.912 ± 0.005	0.910 ± 0.003	$0.855 {\pm} 0.004$	$0.852 {\pm} 0.004$	$0.850 {\pm} 0.006$	0.860 ± 0.004	0.740 ± 0.006	
BioT5	0.937 ± 0.001	0.937 ± 0.004	0.874 ± 0.001	0.870 ± 0.001	0.865 ± 0.002	0.880 ± 0.003	0.765 ± 0.004	
FusionDTI-BAN	$0.923 {\pm} 0.002$	$0.921 {\pm} 0.002$	$0.856 {\pm} 0.001$	$0.857 {\pm} 0.001$	$0.854{\pm}0.002$	$0.858 {\pm} 0.002$	$0.724{\pm}0.001$	
FusionDTI-CAN	$0.951{\pm}0.002$	$0.951{\pm}0.002$	$0.889{\pm}0.002$	$0.890{\pm}0.002$	$0.888{\pm}0.003$	$0.891{\pm}0.002$	$0.778{\pm}0.002$	
Human								
SVM	$0.940{\pm}0.006$	$0.920 {\pm} 0.009$	$0.895 {\pm} 0.010$	$0.892{\pm}0.011$	$0.880{\pm}0.015$	$0.910{\pm}0.009$	$0.800{\pm}0.012$	
RF	$0.952{\pm}0.011$	$0.953{\pm}0.010$	$0.920{\pm}0.012$	$0.915 {\pm} 0.013$	$0.910{\pm}0.017$	$0.930{\pm}0.014$	$0.820 {\pm} 0.009$	
DeepConv-DTI	$0.980{\pm}0.002$	$0.981{\pm}0.002$	$0.927 {\pm} 0.007$	$0.923 {\pm} 0.006$	$0.920 {\pm} 0.009$	$0.930 {\pm} 0.008$	$0.860 {\pm} 0.010$	
GraphDTA	$0.981{\pm}0.001$	$0.982{\pm}0.002$	$0.930{\pm}0.008$	$0.925 {\pm} 0.008$	$0.920 {\pm} 0.011$	$0.935 {\pm} 0.009$	$0.870 {\pm} 0.009$	
MolTrans	$0.980 {\pm} 0.002$	$0.978 {\pm} 0.003$	0.925 ± 0.011	0.920 ± 0.012	0.915 ± 0.016	0.930 ± 0.013	0.855 ± 0.010	
DrugBAN	0.982 ± 0.002	$0.980 {\pm} 0.003$	$0.930 {\pm} 0.004$	0.903 ± 0.003	0.900 ± 0.005	0.908 ± 0.004	0.810 ± 0.004	
SiamDTI	0.970 ± 0.002	0.969 ± 0.003	0.920 ± 0.006	0.915 ± 0.006	0.910 ± 0.008	0.925 ± 0.007	0.840 ± 0.009	
BioT5	0.989 ± 0.001	$\underline{0.985{\pm}0.002}$	$\underline{0.939 \pm 0.008}$	0.937 ± 0.004	$\underline{0.929 \pm 0.010}$	0.941 ± 0.004	0.892 ± 0.006	
FusionDTI-BAN	$0.984{\pm}0.002$	$0.984{\pm}0.003$	$0.938 {\pm} 0.003$	$0.934{\pm}0.002$	$0.927 {\pm} 0.004$	0.931 ± 0.003	0.870 ± 0.003	
FusionDTI-CAN	0.991±0.002	0.989±0.002	0.947±0.002	0.948±0.002	0.955±0.033	0.950±0.031	0.905±0.045	

Table 11: In-domain performance comparison of FusionDTI and the baselines on the BindingDB, Human and BioSNAP datasets (**Best**, <u>Second Best</u>).

Model	AUC	AUPR	Accuracy	F1	Sensitivity	Specificity	MCC	
BindingDB								
SVM	$0.490{\pm}0.015$	$0.460{\pm}0.001$	0.531±0.009	0.521 ± 0.010	$0.508{\pm}0.015$	$0.548{\pm}0.011$	$0.150{\pm}0.012$	
RF	$0.493 {\pm} 0.021$	$0.468 {\pm} 0.023$	$0.535 {\pm} 0.012$	$0.525 {\pm} 0.013$	$0.512{\pm}0.020$	$0.550{\pm}0.014$	$0.162{\pm}0.015$	
GraphDTA	$0.536{\pm}0.015$	$0.496 {\pm} 0.029$	$0.472 {\pm} 0.009$	$0.462{\pm}0.008$	$0.460 {\pm} 0.014$	$0.478 {\pm} 0.011$	$0.100{\pm}0.012$	
DeepConv-DTI	$0.527 {\pm} 0.038$	$0.499 {\pm} 0.035$	$0.490 {\pm} 0.027$	$0.480{\pm}0.026$	$0.475 {\pm} 0.030$	$0.495 {\pm} 0.023$	$0.115 {\pm} 0.020$	
MolTrans	$0.554 {\pm} 0.024$	$0.511 {\pm} 0.025$	$0.470 {\pm} 0.004$	$0.460 {\pm} 0.005$	$0.455{\pm}0.008$	$0.478 {\pm} 0.007$	$0.105 {\pm} 0.008$	
DrugBAN	0.604 ± 0.027	0.570 ± 0.047	0.509 ± 0.021	0.582 ± 0.030	0.565 ± 0.022	0.580 ± 0.025	0.187 ± 0.031	
SiamDTI	$0.627 {\pm} 0.027$	$0.571 {\pm} 0.024$	$0.563 {\pm} 0.033$	$0.550{\pm}0.032$	$0.540 {\pm} 0.036$	$0.580{\pm}0.028$	$0.190 {\pm} 0.030$	
BioT5	$0.651 {\pm} 0.002$	$0.653 {\pm} 0.003$	$0.621 {\pm} 0.005$	$0.608 {\pm} 0.004$	$0.600 {\pm} 0.006$	$0.635 {\pm} 0.005$	0.220 ± 0.007	
FusionDTI-BAN	$0.659 {\pm} 0.002$	$0.663 {\pm} 0.002$	$0.633 {\pm} 0.003$	$0.587 {\pm} 0.002$	$0.603 {\pm} 0.003$	$0.589{\pm}0.002$	$0.276 {\pm} 0.003$	
FusionDTI-CAN	$0.681{\pm}0.005$	$0.680{\pm}0.012$	$0.652{\pm}0.005$	$0.601{\pm}0.005$	$0.628{\pm}0.006$	$0.692{\pm}0.005$	$0.302{\pm}0.005$	
			BioSN	AP				
SVM	$0.602{\pm}0.005$	$0.528 {\pm} 0.005$	$0.513{\pm}0.011$	$0.502{\pm}0.012$	$0.490{\pm}0.014$	$0.523{\pm}0.013$	$0.150{\pm}0.010$	
RF	$0.590{\pm}0.015$	$0.568{\pm}0.018$	$0.499 {\pm} 0.004$	$0.488 {\pm} 0.005$	$0.478 {\pm} 0.008$	$0.513 {\pm} 0.007$	$0.135 {\pm} 0.008$	
GraphDTA	$0.618{\pm}0.005$	$0.618{\pm}0.008$	$0.535 {\pm} 0.024$	$0.528{\pm}0.023$	$0.520{\pm}0.027$	$0.550{\pm}0.020$	$0.170 {\pm} 0.025$	
DeepConv-DTI	$0.645 {\pm} 0.022$	$0.642 {\pm} 0.032$	$0.558 {\pm} 0.025$	$0.550{\pm}0.024$	$0.543 {\pm} 0.030$	$0.573 {\pm} 0.027$	$0.200{\pm}0.028$	
MolTrans	$0.621 {\pm} 0.015$	$0.608 {\pm} 0.022$	$0.546 {\pm} 0.032$	$0.538{\pm}0.031$	$0.530 {\pm} 0.035$	$0.563 {\pm} 0.033$	$0.185{\pm}0.034$	
DrugBAN	$\underline{0.685{\pm}0.004}$	$\underline{0.713{\pm}0.005}$	$\underline{0.692{\pm}0.006}$	$\underline{0.587{\pm}0.005}$	$\underline{0.522{\pm}0.011}$	$\underline{0.690{\pm}0.012}$	$\underline{0.219{\pm}0.017}$	
SiamDTI	$0.718 {\pm} 0.005$	$0.725 {\pm} 0.005$	$0.623 {\pm} 0.007$	$0.610 {\pm} 0.006$	$0.600 {\pm} 0.007$	$0.675 {\pm} 0.006$	$0.240 {\pm} 0.008$	
BioT5	$0.720{\pm}0.008$	$0.718 {\pm} 0.004$	$0.715 {\pm} 0.009$	$0.590{\pm}0.010$	$0.510{\pm}0.012$	$0.710{\pm}0.010$	$0.250 {\pm} 0.011$	
FusionDTI-BAN	$0.723 {\pm} 0.002$	$0.721 {\pm} 0.002$	$0.726 {\pm} 0.001$	$0.597 {\pm} 0.001$	$0.504{\pm}0.012$	0.713 ± 0.011	0.254 ± 0.010	
FusionDTI-CAN	$\textbf{0.748}{\pm 0.021}$	$\textbf{0.766}{\pm}\textbf{0.017}$	$0.734{\pm}0.012$	$0.602{\pm}0.012$	$0.531{\pm}0.013$	$0.736{\pm}0.012$	$\textbf{0.268}{\pm 0.011}$	
Human								
SVM	$0.621 {\pm} 0.036$	$0.637 {\pm} 0.009$	$0.533 {\pm} 0.011$	$0.525{\pm}0.012$	$0.520{\pm}0.015$	$0.546{\pm}0.010$	$0.175 {\pm} 0.011$	
RF	$0.642{\pm}0.011$	$0.663 {\pm} 0.050$	$0.543 {\pm} 0.014$	$0.535 {\pm} 0.015$	$0.530{\pm}0.018$	$0.556{\pm}0.013$	$0.184{\pm}0.012$	
GraphDTA	$0.822 {\pm} 0.009$	$0.759 {\pm} 0.006$	$0.709 {\pm} 0.016$	$0.705 {\pm} 0.017$	$0.702 {\pm} 0.020$	$0.713 {\pm} 0.015$	$0.198 {\pm} 0.017$	
DeepConv-DTI	$0.761 {\pm} 0.016$	$0.628{\pm}0.022$	$0.711 {\pm} 0.030$	$0.704{\pm}0.031$	$0.704{\pm}0.035$	$0.728 {\pm} 0.027$	$0.203 {\pm} 0.030$	
MolTrans	$0.810{\pm}0.021$	$0.745 {\pm} 0.034$	$0.713 {\pm} 0.032$	$0.725 {\pm} 0.033$	$0.720{\pm}0.037$	$0.740{\pm}0.031$	$0.215 {\pm} 0.032$	
DrugBAN	$0.833 {\pm} 0.020$	$0.760 {\pm} 0.031$	$0.709 {\pm} 0.005$	$0.713 {\pm} 0.030$	$0.706 {\pm} 0.022$	$0.720{\pm}0.015$	0.242 ± 0.010	
SiamDTI	$0.863{\pm}0.019$	0.807 ± 0.040	$0.720{\pm}0.010$	$0.729 {\pm} 0.015$	$0.712{\pm}0.020$	$0.736{\pm}0.013$	$0.250 {\pm} 0.015$	
BioT5	$\underline{0.856{\pm}0.003}$	$0.853{\pm}0.003$	$0.715 {\pm} 0.002$	$0.741 {\pm} 0.010$	$0.738{\pm}0.009$	$0.739{\pm}0.013$	$\underline{0.258{\pm}0.013}$	
FusionDTI-BAN	$0.784{\pm}0.002$	$0.790 {\pm} 0.003$	$0.733 {\pm} 0.003$	$0.725 {\pm} 0.002$	$0.713 {\pm} 0.004$	$0.698 {\pm} 0.013$	0.212 ± 0.011	
FusionDTI-CAN	$0.801 {\pm} 0.037$	$0.803 {\pm} 0.032$	$0.738{\pm}0.002$	$\underline{0.736{\pm}0.010}$	$\underline{0.732{\pm}0.013}$	$0.737 {\pm} 0.010$	$0.261{\pm}0.010$	

Table 12: Cross-domain performance comparison of FusionDTI and the baselines on the BindingDB, Human and BioSNAP datasets (**Best**, <u>Second Best</u>).

Method	AUROC	AUPRC	Sensitivity	Specificity
DeepDTA	0.892 ± 0.0066	0.378 ± 0.0231	0.854 ± 0.0066	0.792 ± 0.0291
MolTrans	0.898 ± 0.0050	0.371 ± 0.0067	0.865 ± 0.0050	0.783 ± 0.0387
ML-DTI	0.910 ± 0.0034	0.381 ± 0.0247	0.895 ± 0.0034	0.795 ± 0.0183
DGraphGTA (Alphafold2)	0.885 ± 0.0099	0.316 ± 0.0447	0.894 ± 0.0034	0.724 ± 0.0467
iNGNN-DTI	0.931 ± 0.0027	0.473 ± 0.0167	0.922 ± 0.0155	0.802 ± 0.0240
LANTERN	$\textbf{0.995} \pm \textbf{0.0037}$	0.905 ± 0.0238	$\underline{0.976\pm0.0159}$	$\underline{0.964 \pm 0.0207}$
FusionDTI-BAN	0.973 ± 0.0045	$\underline{0.969 \pm 0.0121}$	0.962 ± 0.0122	0.952 ± 0.0134
FusionDTI-CAN	$\underline{0.987 \pm 0.0032}$	$\textbf{0.978} \pm \textbf{0.0103}$	$\textbf{0.979} \pm \textbf{0.0102}$	$\textbf{0.972} \pm \textbf{0.0116}$

Table 13: Performance comparison on the DAVIS dataset (Best, Second Best).

Method	AUROC	AUPRC	Sensitivity	Specificity
DeepDTA	0.912 ± 0.0037	0.743 ± 0.0127	0.881 ± 0.0056	0.780 ± 0.0127
MolTrans	0.899 ± 0.0022	0.691 ± 0.0142	0.872 ± 0.0116	0.760 ± 0.0160
ML-DTI	0.909 ± 0.0020	0.727 ± 0.0108	0.878 ± 0.0111	0.779 ± 0.0113
DGraphGTA (Alphafold2)	0.911 ± 0.0004	0.739 ± 0.0043	0.881 ± 0.0183	0.784 ± 0.0277
iNGNN-DTI	0.915 ± 0.0016	0.753 ± 0.0071	0.888 ± 0.0183	0.779 ± 0.0146
LANTERN	$\underline{0.976 \pm 0.0154}$	$\underline{0.977 \pm 0.0088}$	$\underline{0.959 \pm 0.0268}$	$\underline{0.965\pm0.0074}$
FusionDTI-BAN	0.974 ± 0.0081	0.976 ± 0.0054	0.952 ± 0.0162	0.947 ± 0.0138
FusionDTI-CAN	$\textbf{0.981} \pm \textbf{0.0064}$	$\textbf{0.981} \pm \textbf{0.0045}$	$\textbf{0.969} \pm \textbf{0.0124}$	$\textbf{0.967} \pm \textbf{0.0156}$

Table 14: Performance comparison on the KIBA dataset (Best, Second Best).

Model	AUC	0.5% RE	1% RE	2% RE	5% RE
DrugVQA	0.972 ± 0.003	88.170 ± 4.88	58.710 ± 2.74	35.060 ± 1.91	17.390 ± 0.94
DrugClip	0.966	118.10	67.17	37.17	16.59
HyperPCM	0.982 ± 0.006	183.04 ± 4.53	91.28 ± 3.35	45.62 ± 2.15	17.13 ± 1.17
MIN	$\underline{0.983 \pm 0.002}$	$\textbf{197.741} \pm \textbf{4.73}$	$\textbf{99.563} \pm \textbf{2.49}$	$\underline{49.926 \pm 1.87}$	$\underline{19.965\pm0.91}$
FusionDTI-BAN	0.9769 ± 0.015	176.8525 ± 2.71	89.2656 ± 2.36	45.9098 ± 1.38	18.5168 ± 0.33
FusionDTI-CAN	$\textbf{0.986} \pm 0.012$	$\underline{186.7469}\pm 6.26$	$\underline{97.8801} \pm 3.50$	$\textbf{52.6352} \pm 2.05$	$\textbf{21.5439} \pm 0.26$

Table 15: Performance comparison on the DUD-E dataset (Best, Second Best).