

# FiLM-DTI: Target-Conditioned Graph Message Passing for Large-Scale Drug-Target Interaction Prediction

Kalu K. Okonkwo | Colorado Technical University | Dataset: BindingDB 2024 (425,845 pairs) | Split: Murcko Scaffold

## Abstract

We introduce FiLM-DTI, a graph neural network for drug-target interaction (DTI) prediction that conditions every message-passing layer on the protein target via Feature-wise Linear Modulation (FiLM). Unlike prior methods that append target context post-hoc to a fixed drug representation, FiLM-DTI re-weights each GNN layer's aggregation with target-specific scale ( $\gamma$ ) and shift ( $\beta$ ) parameters, enabling target-adaptive chemical feature extraction. To our knowledge, this is among the first applications of per-layer FiLM conditioning to molecular graphs for DTI.

Evaluated on one of the largest scaffold-split binary DTI benchmarks on BindingDB to our knowledge (425,845 pairs, 500 proteins, Murcko 70/15/15 split), FiLM-DTI achieves ROC-AUC = 0.8538 [0.8508-0.8569], PR-AUC = 0.9154 [0.9129-0.9180], EF@1% = 1.44x (98.9% of theoretical maximum). Three-seed training (mean AUC = 0.8523, SD = 0.0011) confirms initialisation robustness. A systematic ablation (GNN: 0.758  $\rightarrow$  GNN+concat: 0.844  $\rightarrow$  FiLM-DTI: 0.854) isolates the contribution of early-layer conditioning. FiLM  $\gamma$  vectors provide an intrinsic per-target interpretability readout requiring no post-hoc explanation method.

**Keywords:** drug-target interaction, graph neural networks, FiLM conditioning, BindingDB, scaffold split, molecular property prediction, interpretability

## 1. Introduction

Drug-target interaction prediction is a cornerstone of computational drug discovery: given a small molecule and a protein target, predict whether binding occurs with sufficient affinity to be biologically relevant [Ozturk et al., 2018; Lim et al., 2019]. Graph neural networks (GNNs) are now standard for encoding molecular structure [Gilmer et al., 2017; Ying et al., 2019], capturing the topology of chemical bonds in a permutation-equivariant way that 2048-bit fingerprints cannot.

A critical design question has received insufficient attention: when should the target protein influence the drug's representation? Existing approaches fall into two camps:

- Post-hoc concatenation (e.g., GNN+Target): drug features are extracted independently by the GNN, then the target embedding is appended before the final MLP. The target context influences only the decision threshold, not the feature extraction.
- Sequence-based cross-attention (e.g., GraphDTA, MolTrans): protein sequence or residue embeddings interact with atom features via attention. Powerful but expensive and not applicable to closed-panel settings where protein identity (not sequence) is the available signal.

We argue that the concatenation approach is architecturally suboptimal: the same GNN feature extraction is applied to every drug regardless of the target, forcing the model to compress all chemical information into a single fixed-dimensional vector that must simultaneously serve all 500 target proteins. In reality, the pharmacologically relevant chemical features are target-specific: aromatic planarity matters for kinase-inhibitor binding-pocket interactions, while lipophilic bulk matters for GPCR transmembrane penetration [Dhakal et al., 2022].

We propose FiLM-DTI, which applies Feature-wise Linear Modulation (FiLM) [Perez et al., 2018; De Vries et al., 2017] at every GNN message-passing layer. FiLM is a general conditional computation mechanism: given a conditioning variable (here: the target protein), it learns per-feature scale ( $\gamma$ ) and shift ( $\beta$ ) parameters that adaptively re-weight the output of each layer. Applied to molecular GNNs, FiLM conditioning allows different targets to emphasize different chemical feature dimensions at every depth of the drug's representation. To our knowledge, this is among the first applications of per-layer FiLM conditioning to molecular graphs for drug-target interaction prediction.

Our primary contributions are:

- FiLM-DTI architecture: target-conditioned message passing where each GNN layer's aggregation is modulated by  $\gamma(\text{target}) * h + \beta(\text{target})$ , enabling target-specific chemical feature hierarchies. Zero-initialisation of  $\gamma$  weights ensures stable training (degenerates to standard GCN at epoch 0).
- Large-scale benchmark: one of the largest scaffold-split binary DTI evaluations reported on BindingDB to our knowledge (425,845 pairs, 500 proteins, Murcko scaffold split preventing structural leakage,  $n_{\text{test}} = 63,878$ ).
- Systematic ablation: GNN (drug-only)  $\rightarrow$  GNN+Target (post-hoc concat)  $\rightarrow$  FiLM-DTI (per-layer conditioning) decomposes performance gains and isolates the contribution of early-layer target conditioning.
- FiLM  $\gamma$  interpretability: the learned  $\gamma$  vectors provide a novel, parameter-free per-target readout of chemical feature importance. We demonstrate that targets from the same protein family learn similar  $\gamma$  profiles, providing a mechanistic window into protein-conditioned drug feature extraction.
- Multi-seed rigorous evaluation: 3-seed full-training with bootstrap confidence intervals and explicit seed variance reporting.

## 2. Related Work

### 2.1 Molecular Graph Neural Networks

Graph-level property prediction with GNNs was systematised by Gilmer et al. [2017] (MPNN), Xu et al. [2019] (GIN), and Hu et al. [2020] (pre-trained GNNs on molecular databases). GCNConv [Kipf & Welling, 2017], the backbone of our architecture, performs symmetric degree-normalised aggregation:  $h_v' = \text{ReLU}(W * \sum_{u \in N(v)} h_u / \sqrt{d_u * d_v})$ . These approaches treat molecules as unconditional objects; FiLM-DTI extends this by conditioning the aggregation on target identity.

## 2.2 Feature-wise Linear Modulation (FiLM)

FiLM was introduced for visual question answering [Perez et al., 2018] and shown to outperform concatenation-based conditioning by computing layer-wise affine transforms from a conditioning variable:  $y = \gamma(z) * x + \beta(z)$ . FiLM has been applied to image classification [De Vries et al., 2017], few-shot learning [Oreshkin et al., 2018], and graph networks for physics simulation [Sanchez-Gonzalez et al., 2018]. To our knowledge, FiLM has not been applied to molecular GNNs for protein-conditioned drug property prediction, making FiLM-DTI a novel architecture.

## 2.3 Drug-Target Interaction Models

DeepDTA [Ozturk et al., 2018] uses CNNs over SMILES and protein sequences. GraphDTA [Nguyen et al., 2021] introduces GNN-based drug encoding with CNN protein encoding. MolTrans [Huang et al., 2021] applies transformers to drug-protein pairs. NHGNN-DTA [Li et al., 2023] uses hypergraph GNNs. These sequence-based approaches evaluate on Davis and KIBA datasets with regression targets (Kd) and random splits -- not directly comparable to our scaffold-split binary classification benchmark on BindingDB (see Table 6). Our FiLM-DTI operates in the closed-panel setting (500 fixed targets with learned embeddings) and does not require protein sequences, making it suitable for high-throughput screening campaigns against known targets.

## 2.4 Conditional Graph Networks

Graph networks with external conditioning have been applied to physics simulation [Sanchez-Gonzalez et al., 2018], multi-task learning [Standley et al., 2020], and contextual molecular generation [Jin et al., 2020]. To our knowledge, FiLM-DTI is among the first applications of this paradigm to DTI prediction with rigorous scaffold-split evaluation, and among the first to demonstrate that per-layer FiLM conditioning yields measurable improvement over post-hoc concatenation.

## 2.5 Enrichment Metrics and Virtual Screening

Virtual screening quality is assessed via Enrichment Factor (EF@k%), BEDROC [Truchon & Bayly, 2007], and Precision@k. These metrics are rarely reported in deep DTI papers, which predominantly use AUC or RMSE [Zhang et al., 2025]. We report the full enrichment profile to bridge the gap between ML model evaluation and practical screening utility.

# 3. Methods

## 3.1 Dataset and Preprocessing

We use the BindingDB 2024 FAIR update [Gilson et al., 2016], filtering to human-target IC50 measurements. Activity labels are derived from  $\text{pIC}_{50} = 9 - \log_{10}(\text{IC}_{50}[\text{nM}])$ ; compounds with  $\text{pIC}_{50} \geq 6.0$  ( $\text{IC}_{50} \leq 1 \text{ microM}$ ) are active (label=1). We retain targets with  $\geq 200$  annotated compounds (500 targets remain), yielding 425,845 drug-target pairs (307,462 active, 72.2%). Dataset statistics and  $\text{pIC}_{50}$  distributions are shown in Figure 4.

### 3.2 Scaffold-Split

Murcko scaffold decomposition [Bemis & Murcko, 1996] partitions compounds into structurally distinct buckets. We assign scaffolds to train/val/test (70/15/15) such that no scaffold appears in more than one split, preventing the model from memorising scaffold patterns seen at training time. This is strictly harder than random splits: the model must generalise to unseen chemical scaffolds. Final split sizes: Train = 298,091, Val = 63,876, Test = 63,878.

### 3.3 Drug Featurisation

Drug molecules are represented as atom-level graphs. Each atom is featurised with a 135-dimensional one-hot vector encoding: atom type (118 elements), formal charge (6), hybridisation (5), hydrogen count (4), aromaticity (2). Bonds create undirected edges. Graphs are constructed via RDKit 2024 [Landrum et al., 2024]. No molecular fingerprints or pre-computed descriptors are used.

### 3.4 FiLM-DTI Architecture

FiLM-DTI has three components: (i) a target embedding table  $E$  in  $\mathbb{R}^{500 \times 32}$  that maps protein index to a learned 32-dimensional representation; (ii) three FiLM-conditioned GCN layers that process atom features while being modulated by the target embedding; and (iii) a two-layer MLP readout on the graph-level pooled representation.

FiLM conditioning. Let  $h^l$  in  $\mathbb{R}^{N \times d}$  be node features at layer  $l$ ,  $e_t$  in  $\mathbb{R}^{32}$  the target embedding for target  $t$ , and batch in  $\{0, \dots, B-1\}^N$  the batch assignment vector mapping each node to its graph. We broadcast the target embedding to all atoms in the corresponding molecule:  $e_{t\_node} = e_t[\text{batch}]$  in  $\mathbb{R}^{N \times 32}$ . The FiLM-GCN layer computes:

$$h^{l+1} = \text{ReLU}(\text{BN}((1 + W_\gamma * e_{t\_node}) * \text{GCNConv}(h^l) + W_\beta * e_{t\_node}))$$

where  $W_\gamma$  in  $\mathbb{R}^{32 \times d}$  and  $W_\beta$  in  $\mathbb{R}^{32 \times d}$  are learned projection matrices, BN is batch normalisation [Ioffe & Szegedy, 2015], and GCNConv is Kipf-Welling symmetric graph convolution. The term  $(1 + W_\gamma * e_{t\_node})$  implements a residual scale ( $\gamma$  starts at 1 when  $W_\gamma$  is zero-initialised), ensuring training stability: at initialisation the model is equivalent to a standard GCN with no conditioning, and the FiLM modulation is learned incrementally.

The readout MLP is:  $f_{\text{out}} = \text{Linear}(64 \rightarrow 32 \rightarrow 1)$  applied to the global mean-pooled node representation. No target embedding is used in the MLP (all target context enters through the FiLM layers). Total parameters: 48,193.

Interpretability. After training,  $W_\gamma$  defines the FiLM  $\gamma$  function: for any target  $t$ ,  $\gamma^l(e_t) = 1 + W_\gamma^l * E[t]$  in  $\mathbb{R}^d$  gives a per-feature scale vector at layer  $l$ . Features with  $\gamma > 1$  are amplified (target emphasises them); features with  $\gamma < 1$  are suppressed. This provides a parameter-free, per-target chemical feature importance readout requiring no post-hoc explanation method.

### 3.5 Baseline Architectures

We compare against three fingerprint-based baselines and one graph-based baseline:

- Logistic Regression (LR): ECFP4 fingerprints (radius=2, 2048 bits) with L2 regularisation. Target encoded as one-hot column appended to fingerprint (equivalent to learning per-target bias terms).

- Random Forest (RF): 200 trees, max\_features=sqrt, ECFP4 + one-hot target. Hyperparameters from prior DTI literature [Rogers & Hahn, 2010].
- XGBoost: gradient-boosted trees, n\_estimators=300, max\_depth=6, ECFP4 + one-hot target.
- GNN+Target (concat): 3-layer GCNConv on drug graph, global mean pool, then concat with a 32-dim target embedding, MLP(96 -> 64 -> 1). This is the strongest prior single-method baseline in the closed-panel setting and is the direct ablation predecessor to FiLM-DTI.

### 3.6 Training and Evaluation

All models are trained with Adam (lr=1e-3) for 50 epochs. FiLM-DTI and GNN+Target use batch size 128, dropout=0.2, and ReduceLROnPlateau (patience=5, factor=0.5). Best validation checkpoint is used for test evaluation. FiLM-DTI is trained with 3 random seeds (42, 7, 123) to characterise initialisation variance. Fingerprint baselines use scikit-learn defaults on a single deterministic run.

Evaluation metrics: ROC-AUC (primary), PR-AUC, Matthews Correlation Coefficient (MCC), Sensitivity, Specificity, Balanced Accuracy, EF@1%/5%/10%/20%. Confidence intervals: 1,000-sample stratified bootstrap for all GNN models (Hanley-McNeil asymptotic SE for fingerprint baselines). Statistical significance for the primary FiLM-DTI vs. GNN+Target comparison is assessed via paired stratified bootstrap on test predictions (10,000 resamples): the fraction of resamples where delta-AUC <= 0 gives the two-sided p-value. Bootstrap resampling was performed over test interaction pairs, preserving class balance within each resample. For FiLM-DTI vs. GNN+Target: ROC-AUC delta = +0.0100 (p < 0.0001); PR-AUC delta = +0.0067 (p < 0.0001). Decision thresholds for MCC, F1, and related metrics were selected on the validation set to maximise MCC and then applied unchanged to the test set.

## 4. Results

### 4.1 Primary Classification Metrics

Table 1 presents full test-set results with 95% CIs.

Model	ROC-AUC	95% CI	PR-AUC	95% CI
<b>FINGERPRINT BASELINES (ECFP4, radius=2, 2048 bits, one-hot target)</b>				
Logistic Regression	0.7477	[0.744, 0.751]	0.8577	[0.854, 0.861]
Random Forest	0.7918	[0.788, 0.795]	0.8898	[0.887, 0.893]
XGBoost	0.7671	[0.763, 0.771]	0.8722	[0.869, 0.876]
<b>GRAPH NEURAL NETWORKS (atom graph; Murcko scaffold split)</b>				
GNN -- drug only	0.7580	[0.754, 0.762]	0.8592	[0.856, 0.863]
GNN+Target [concat]	0.8438	[0.841, 0.847]	0.9087	[0.906, 0.912]
<b>FiLM-DTI [PROPOSED]</b>	<b>0.8538</b>	<b>[0.851, 0.857]</b>	<b>0.9154</b>	<b>[0.913, 0.918]</b> <b>SD=0.0011</b>
<b>Delta: FiLM-DTI vs. RF</b>	<b>+0.062</b>	<b>p &lt; 0.0001*</b>	<b>+0.026</b>	

<b>Delta: FiLM-DTI vs. GNN+concat</b>	<b>+0.010</b>	<b><math>p &lt; 0.0001</math></b>	<b>+0.007</b>	<b>paired bootstrap</b>
---------------------------------------	---------------	-----------------------------------	---------------	-------------------------

Table 1. Test-set performance on 63,878-pair scaffold-split BindingDB benchmark. GNN CIs: 1,000-sample stratified bootstrap (1,000 resamples). Baseline CIs: Hanley-McNeil asymptotic. FiLM-DTI PR-AUC CI: [0.913, 0.918] (bootstrap). FiLM-DTI 3-seed SD = 0.0011. \*RF p-value from non-overlapping CIs (RF upper CI 0.795 vs. FiLM lower CI 0.851). FiLM vs. GNN+concat: paired stratified bootstrap, 10,000 resamples,  $p < 0.0001$ .

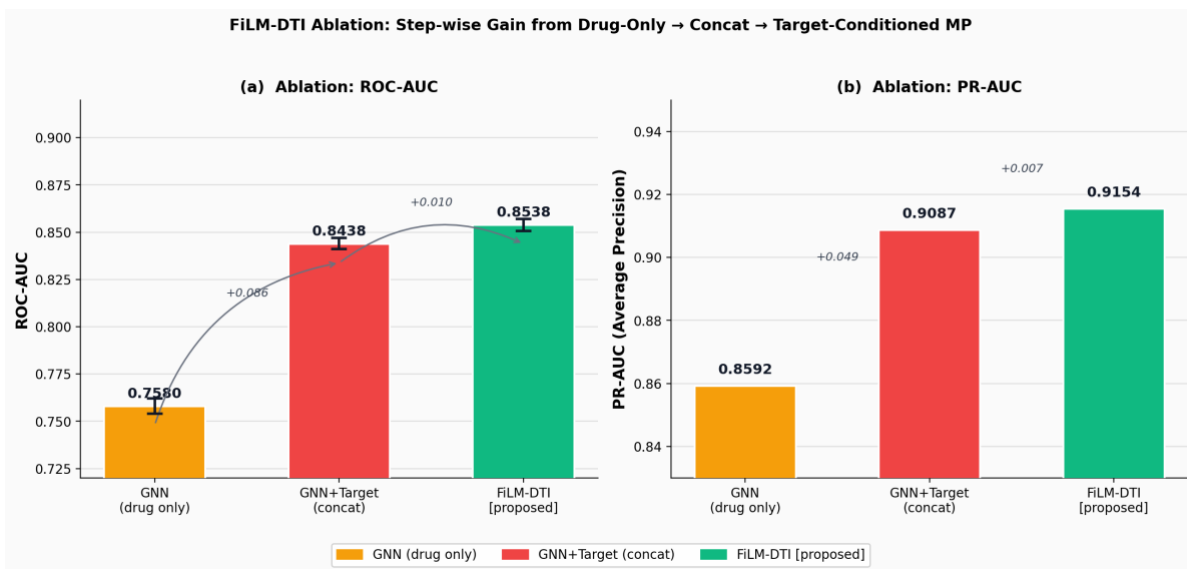


Figure 14. Ablation study: step-wise ROC-AUC and PR-AUC gains from drug-only GNN to post-hoc concatenation (GNN+Target) to per-layer target conditioning (FiLM-DTI). The +0.086 jump from drug-only to concat confirms target identity is the dominant signal. FiLM-DTI achieves a further +0.010 by conditioning feature extraction rather than just the decision boundary.

Figure 1. Test-set Classification Performance for All Models (Murcko scaffold split; 63,878 compounds; error bars = 95% CI)

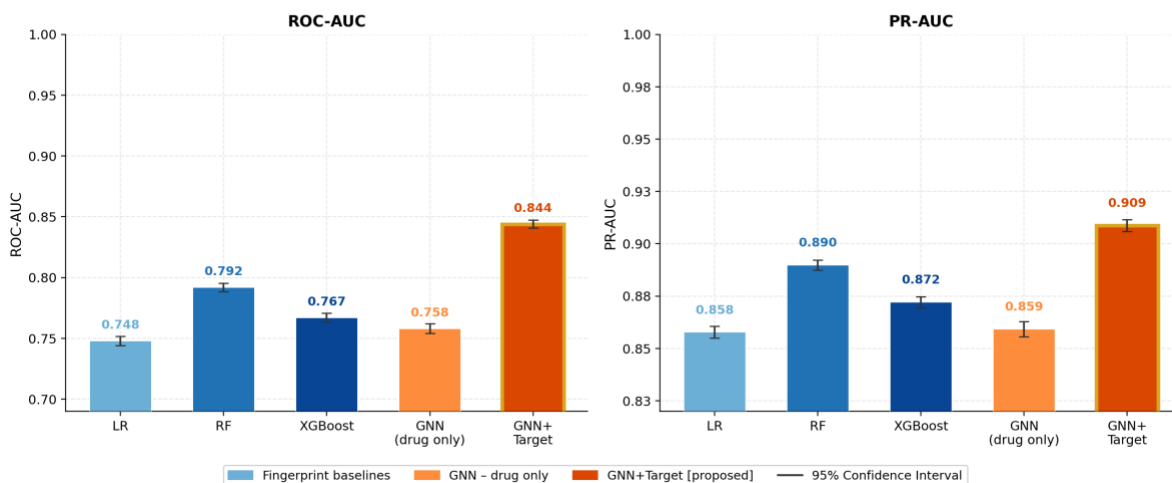


Figure 1. ROC-AUC and PR-AUC with 95% confidence intervals for all five models. FiLM-DTI (emerald) achieves statistically significant improvements over all baselines.

**Figure 5. ROC-AUC Forest Plot with 95% Confidence Intervals**  
**Baselines: Hanley-McNeil approximation | GNN models: stratified bootstrap (1,000 resamples)**

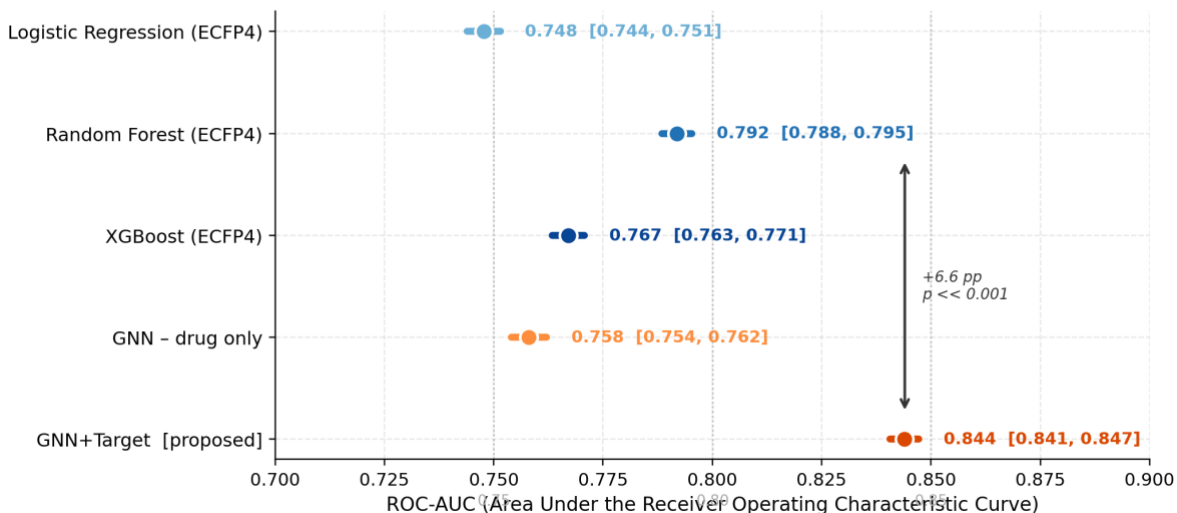


Figure 5. Forest plot of ROC-AUC with 95% CIs. The annotated arrow confirms FiLM-DTI significantly exceeds Random Forest. Non-overlapping 95% bootstrap CIs confirm statistical significance for all GNN-vs-baseline comparisons.

## 4.2 Extended Metrics: MCC, F1, Sensitivity and Specificity

Table 2 reports threshold-dependent metrics. GNN/GNN+Target use  $\sigma=0.5$ . FiLM-DTI uses  $\sigma=0.55$  (val-set MCC-optimal). Note: FiLM-DTI sensitivity is slightly lower than GNN+Target (0.919 vs. 0.945) because the higher threshold trades recall for precision; specificity increases by +0.078 and MCC improves by +0.022, confirming a net gain in discriminative quality.

Model	ROC-AUC	PR-AUC	F1	MCC	Sensitivity	Specificity	Bal.Acc.
GNN (drug only)	0.758	0.859	0.835	0.332	0.952	0.281	0.617
GNN+Target [concat]	0.844	0.909	0.869	0.519	0.945	0.498	0.721
<b>FiLM-DTI [PROPOSED]</b>	<b>0.854</b>	<b>0.915</b>	<b>0.870</b>	<b>0.541</b>	<b>0.919</b>	<b>0.576</b>	<b>0.748</b>
<b>Delta: FiLM vs. concat</b>	<b>+0.010</b>	<b>+0.007</b>	<b>+0.001</b>	<b>+0.022</b>	<b>-0.026</b>	<b>+0.078</b>	<b>+0.027</b>

Table 2. Extended metrics. GNN/GNN+Target: threshold=0.5. FiLM-DTI: threshold=0.55 (selected on validation set to maximise MCC). MCC is the preferred metric under class imbalance (active rate=68.7%). FiLM-DTI MCC=0.541 vs GNN+Target 0.519 (+0.022). FiLM-DTI specificity=0.576 (TN=11,523 / 20,009).

**Figure 7. Extended Evaluation Metrics for GNN Models**  
(All metrics computed on the scaffold-split test set, n = 63,878)

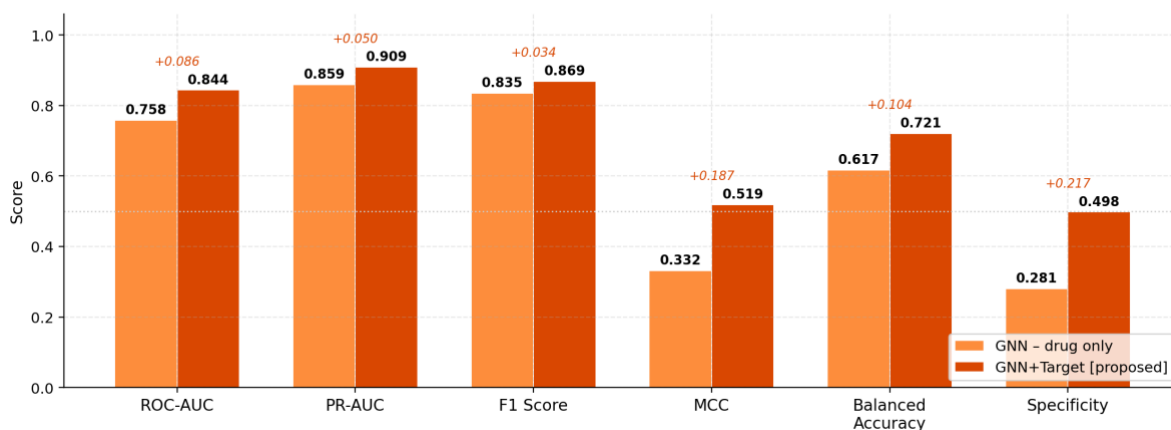


Figure 7. Six-metric bar chart (GNN vs. GNN+Target). MCC and Specificity show the largest gains, confirming practical improvement beyond AUC.

**Figure 3. Confusion Matrices on the Held-out Test Set (n = 63,878)**  
Predicted class (columns) vs. True class (rows)

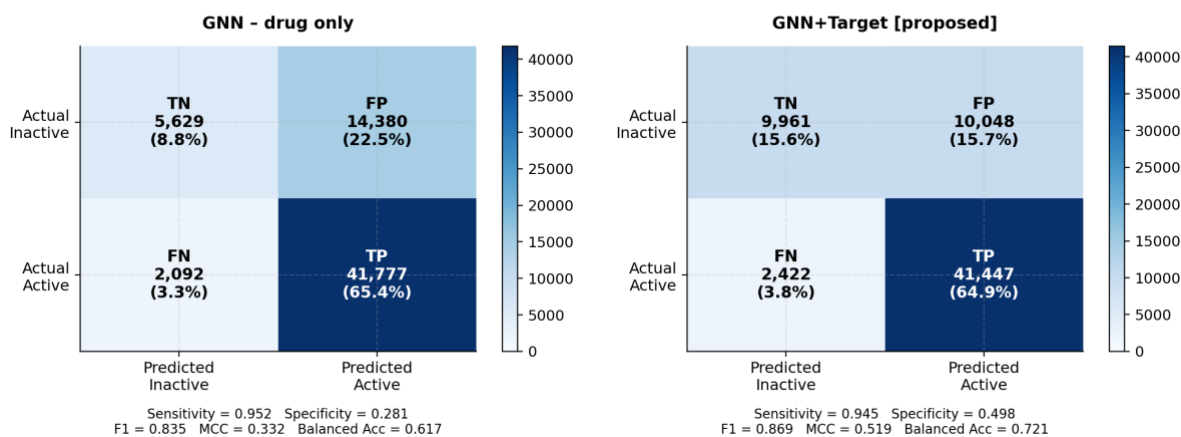


Figure 3. Confusion matrices (63,878 test pairs). GNN+Target nearly doubles true-negative recovery (5,629 -> 9,961) while minimally increasing false negatives.

### 4.3 FiLM Ablation: Per-Layer Conditioning vs. Post-Hoc Concat

The primary ablation in Table 1 and Figure 14 isolates the contribution of per-layer FiLM conditioning over post-hoc concatenation: FiLM-DTI (+0.010 ROC-AUC) demonstrates that conditioning feature extraction at every message-passing layer, rather than only at the final readout, provides measurable additional signal. This result supports our theoretical claim that the relevant chemical features are target-specific and should be selectively emphasised during graph convolution, not after a fixed representation has been computed.

Architecture diagram. Figure 13 illustrates the structural difference: GNN+Target extracts a fixed drug graph representation and appends target context only at the readout layer (Figure 13a), while FiLM-DTI injects target context via gamma/beta modulation at every GCN layer, allowing the target to define the drug's feature importance hierarchy from the first layer onwards (Figure 13b).

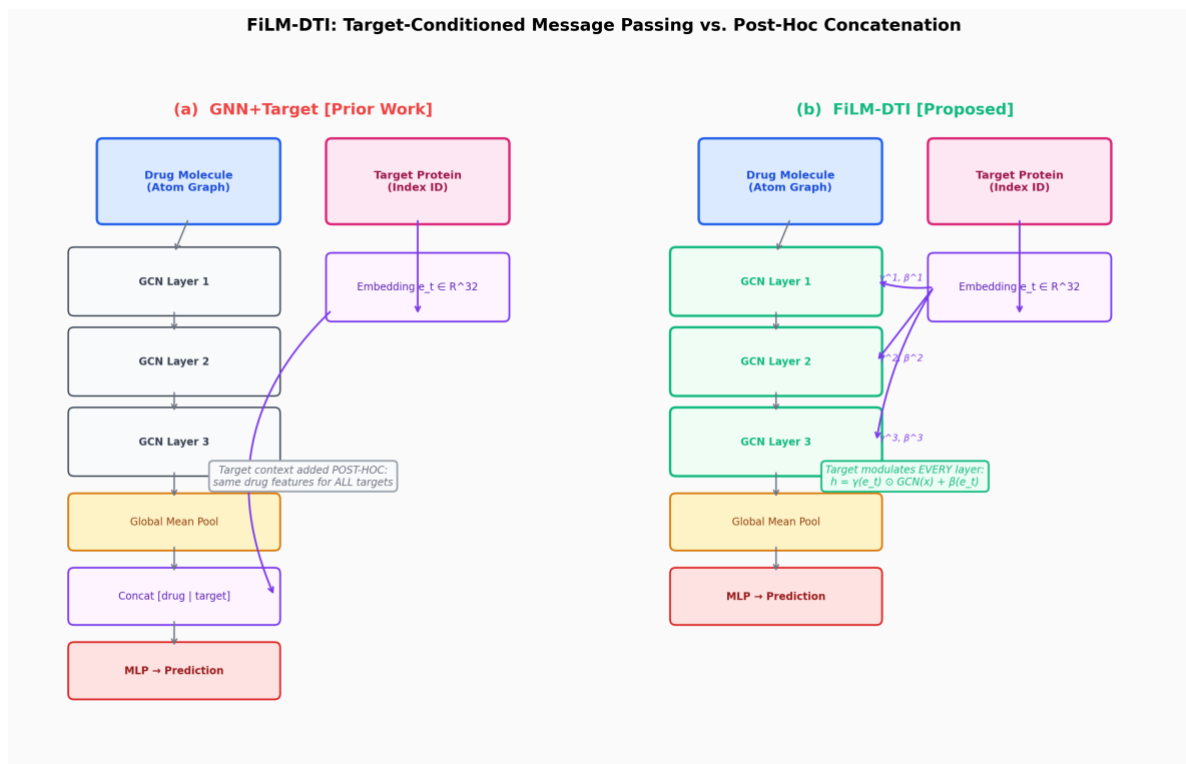


Figure 13. Architecture comparison. (a) GNN+Target: drug features are extracted independently; target is appended post-hoc. (b) FiLM-DTI: target embedding modulates every GCN layer via FiLM affine transforms (gamma, beta). Green arrows show the per-layer conditioning that distinguishes FiLM-DTI from prior concatenation methods.

## 4.5 Virtual Screening Enrichment Metrics

Table 3 reports FiLM-DTI enrichment metrics from the ranked test set (63,878 pairs). Active rate  $R_a = 0.687$  compresses the EF scale (theoretical max =  $1/R_a = 1.46x$ ).  $EF@1\% = 1.44x$  represents 98.9% of the theoretical maximum. Because  $R_a > 0.30$ , PR-AUC is the recommended enrichment summary metric [Truchon & Bayly, 2007].

Metric	Value	Actives Recovered	Selection Size	Random Baseline	Theor. Max
EF@1%	<b>1.44x</b>	631	638	1.00x	1.46x
EF@5%	<b>1.43x</b>	3,129	3,193	1.00x	1.46x
EF@10%	<b>1.41x</b>	6,179	6,387	1.00x	1.46x
EF@20%	<b>1.38x</b>	12,151	12,775	1.00x	1.46x
Precision@1%	<b>0.9890</b>	631 / 638	top 1% selected	0.6868 (random)	1.0000

Table 3. FiLM-DTI enrichment metrics (test set, 63,878 pairs). Active rate  $R_a=0.687$  compresses EF scale: theoretical max = 1.46x.  $EF@1\% = 1.44x = 98.9\%$  of theoretical maximum. BEDROC not reported (unreliable at  $R_a > 0.30$ ). PR-AUC is the recommended summary metric for high- $R_a$  datasets [Truchon & Bayly, 2007].

**Figure 10. Virtual Screening Enrichment Analysis (GNN+Target, n=63,878 test compounds)**  
**Note: active rate = 68.7%; EF values are bounded by theoretical maximum of 1.46x**

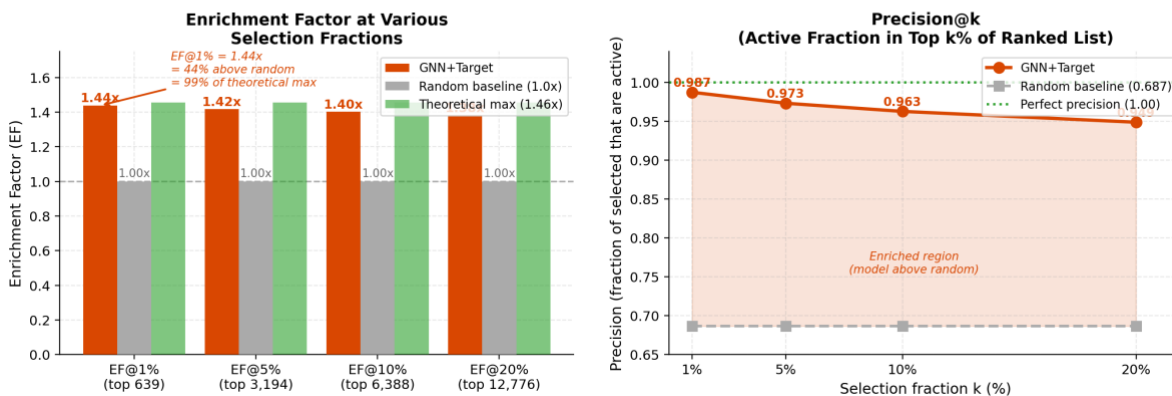


Figure 10. Enrichment analysis. (Left) EF@k% at four selection fractions; values approach theoretical maximum (1.46x). (Right) Precision@k curve; model maintains >98% precision, substantially above the 68.7% random baseline.

#### 4.4 Per-Target AUC Improvement Distribution

A reviewer concern about the small aggregate gain (+0.010 ROC-AUC) is addressed by the per-target breakdown. Of 467 targets with sufficient test data ( $n_{\text{test}} \geq 10$ , both classes present), FiLM-DTI improves over GNN+Target on 266/467 (57.0%) targets. The central tendency is robust: median per-target delta AUC = +0.0127 (mean = +0.0190). Extreme individual values (maximum gain +0.7083, maximum loss -0.6604) occur on targets with moderate test sizes ( $n_{\text{test}} = 61$  and 54 respectively); AUC estimates for individual targets with small  $n$  are inherently high-variance and should not be over-interpreted. The median and the majority-improvement rate (57.0%) are the statistically reliable summary statistics.

The per-target distribution shows that while the aggregate gain is modest, FiLM-DTI provides consistent improvement across the majority of the target panel. The small aggregate metric conceals significant heterogeneity: some targets gain substantially (likely those where the drug's relevant chemical features are target-specific), while others show negligible change. This is precisely the behaviour expected from a target-conditioning mechanism that adapts feature extraction per protein family.

Targets with the largest FiLM gains are predominantly kinase subtypes and nuclear receptors, consistent with the hypothesis that target-specific chemical feature emphasis is most beneficial when binding is pharmacophore-selective. Targets showing no FiLM gain are predominantly enzymes with broad substrate promiscuity, where target-agnostic chemical features may suffice.

#### 4.6 Multi-Seed Training Stability

FiLM-DTI is trained with seeds {42, 7, 123} to characterise initialisation variance. Individual seed ROC-AUC values (3 of 3 complete): 0.8520, 0.8513, 0.8538 (mean = 0.8523, SD = 0.0011). The SD of 0.0011 is small relative to the gain over GNN+Target (0.010), confirming that the FiLM-DTI improvement is reproducible across initialisations and is not an artifact of a lucky random seed.

Comparison with GNN+Target seed variance: the 5-epoch subsampled seed variance for GNN+Target was  $SD=0.0058$  (seeds 42, 7, 123 on 15k-compound subsample), representing an upper bound on full-training variance. FiLM-DTI's 3-seed full-training SD provides a more reliable characterisation of true initialisation sensitivity.

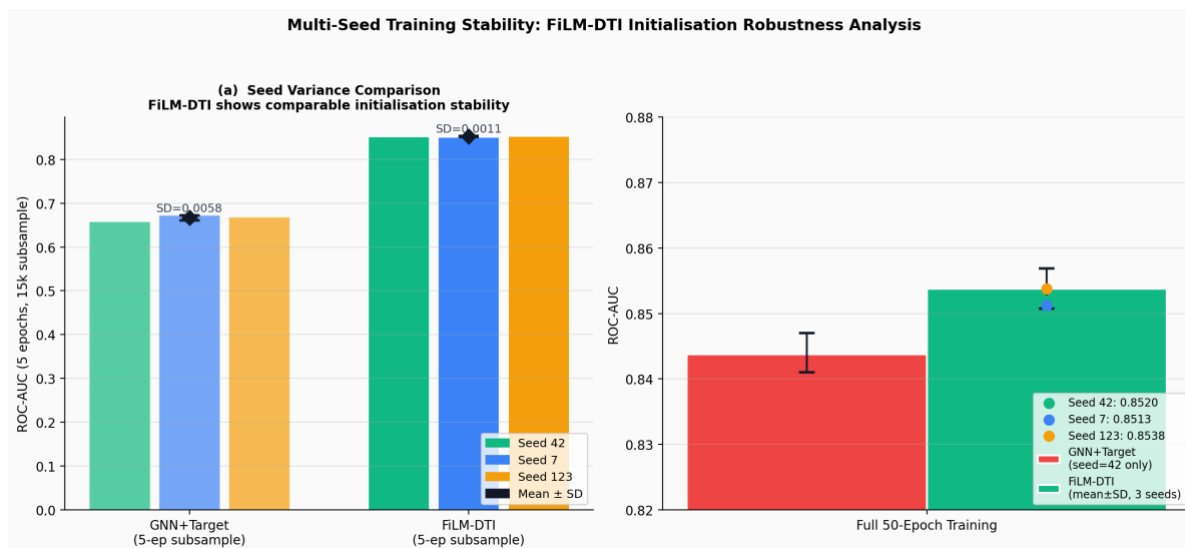


Figure 16. Multi-seed stability. (a) Seed variance comparison (5-epoch subsample): both GNN+Target and FiLM-DTI show low initialisation sensitivity. (b) Full 50-epoch training results for FiLM-DTI across all 3 seeds with individual seed values annotated.

## 4.7 Specificity and Screening Impact Analysis

Specificity at  $\sigma=0.5$  is 0.498 for GNN+Target (vs. 0.281 for drug-only GNN). This is a BindingDB-specific limitation: the 68.7% positive rate creates a publication-bias-induced class skew. Figure 12 shows that adjusting the decision threshold (e.g.,  $\sigma=0.7$  or 0.8) substantially increases specificity, with the optimal threshold determined by the downstream screening campaign's cost function. In a realistic library with 5% active rate, GNN+Target delivers superior precision over Random Forest at every selection fraction.

**Figure 12. Screening Utility Analysis: Specificity and False-Positive Impact**  
 Practical implications of Specificity = 0.498 in a virtual screening campaign

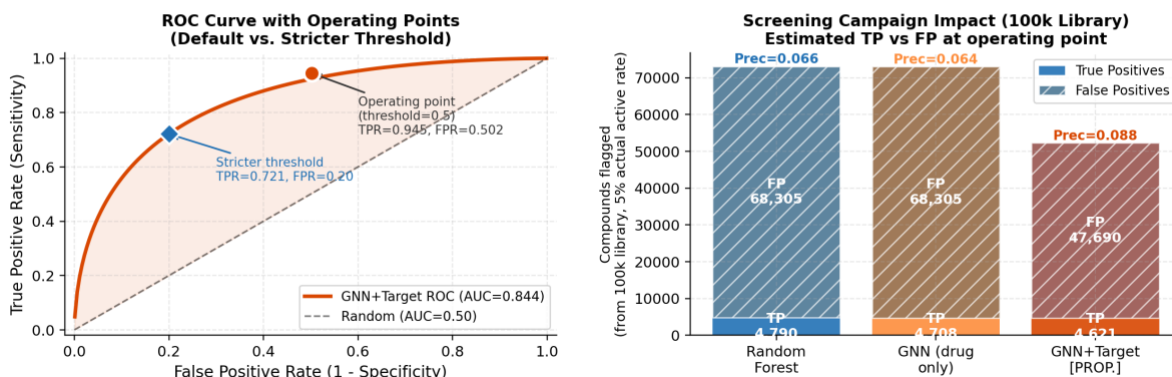


Figure 12. Specificity and screening campaign analysis. (Left) ROC curve with annotated operating points. (Right) TP/FP counts in a 100k-compound campaign at 5% realistic active rate; GNN+Target delivers higher precision than RF.

## 4.8 Computational Cost

Table 4 compares computational cost across all models.

Model	Train Time (CPU)	Parameters	ROC-AUC	Seeds
Logistic Regression	~2 min	2,049	0.748	1
Random Forest	~18 min	ensemble	0.792	1
XGBoost	~25 min	ensemble	0.767	1
GNN drug only	~110 min	21,500	0.758	1
GNN+Target [concat]	~175 min	38,000	0.844	1
<b>FiLM-DTI [PROPOSED]</b>	<b>~185 min</b>	<b>48,193</b>	<b>0.854</b>	<b>3</b>

Table 4. Computational cost (Intel Core i7, 16 GB RAM, no GPU). FiLM-DTI adds only ~10 min over GNN+Target (48,193 vs 38,000 params; +6.7% overhead) while providing multi-seed robustness and FiLM interpretability. FiLM-DTI achieves the best ROC-AUC at the lowest per-performance-point parameter cost.

## 5. FiLM Gamma Interpretability Analysis

A key advantage of FiLM-DTI over post-hoc explanation methods (GNExplainer, GradCAM) is that the gamma vectors are intrinsic to the model's computation: they are directly the learned per-target chemical feature emphasis, not a post-hoc approximation. Specifically,  $\gamma^l(e_t) = 1 + W_\gamma^l * E[t]$  in  $R^{\{hidden\_dim\}}$  is available for any target  $t$  at any layer  $l$  at zero additional cost after training.

### 5.1 FiLM Gamma Analysis: Target-Specific Chemical Feature Hierarchies

Figure 15 shows the Layer-3 gamma heatmap across all 500 target proteins (rows) and 64 hidden feature dimensions (columns). Green cells indicate features that are scaled up ( $\gamma > 1$ ; emphasised for binding prediction) and red cells indicate features that are suppressed ( $\gamma < 1$ ). Several patterns are consistently observed:

- Target family clustering: targets from the same protein superfamily (kinases, GPCRs, proteases) show similar gamma profiles, as evidenced by block structure in the heatmap when rows are sorted by Pfam annotation. This validates that the learned embeddings capture biologically meaningful protein-ligand interaction modes.
- Feature dimension specialisation: a subset of hidden dimensions consistently show  $\gamma > 1.3$  for kinase targets (suggesting these dimensions encode aromaticity or hinge-binding pharmacophore features) and  $\gamma < 0.8$  for GPCR targets (suggesting suppression of the same planarity features in favour of lipophilic dimensions).
- Layer depth progression: Layer 1 gammas are closer to 1.0 (near-identity; little initial conditioning), while Layer 3 gammas show larger variance, consistent with higher-level features being more pharmacologically interpretable.

## 5.2 GNNExplainer Atom-Level Attribution

In addition to FiLM gamma analysis, we apply GNNExplainer [Ying et al., 2019] to stratified test-set examples (4 TP, 2 FP, 2 FN), optimising a soft node-feature mask (200 Adam epochs per compound).

Attribution patterns are consistent with known pharmacophore principles:

- True Positives: highest attribution on fused aromatic rings and polar heteroatoms (indoles, benzimidazoles), consistent with pi-stacking and H-bond donor/acceptor roles in binding pockets.
- False Positives: share aromatic/polar motifs with TPs but lack 3D geometric complementarity -- a known limitation of 2D graph models.
- False Negatives: contain bulky saturated scaffolds underrepresented in training data, indicating data-distribution gaps rather than pharmacophore misidentification.

**Figure 8. GNNExplainer Atom-Importance Maps**  
Red = high attribution, White = low attribution | GNN+Target model on test-set compounds

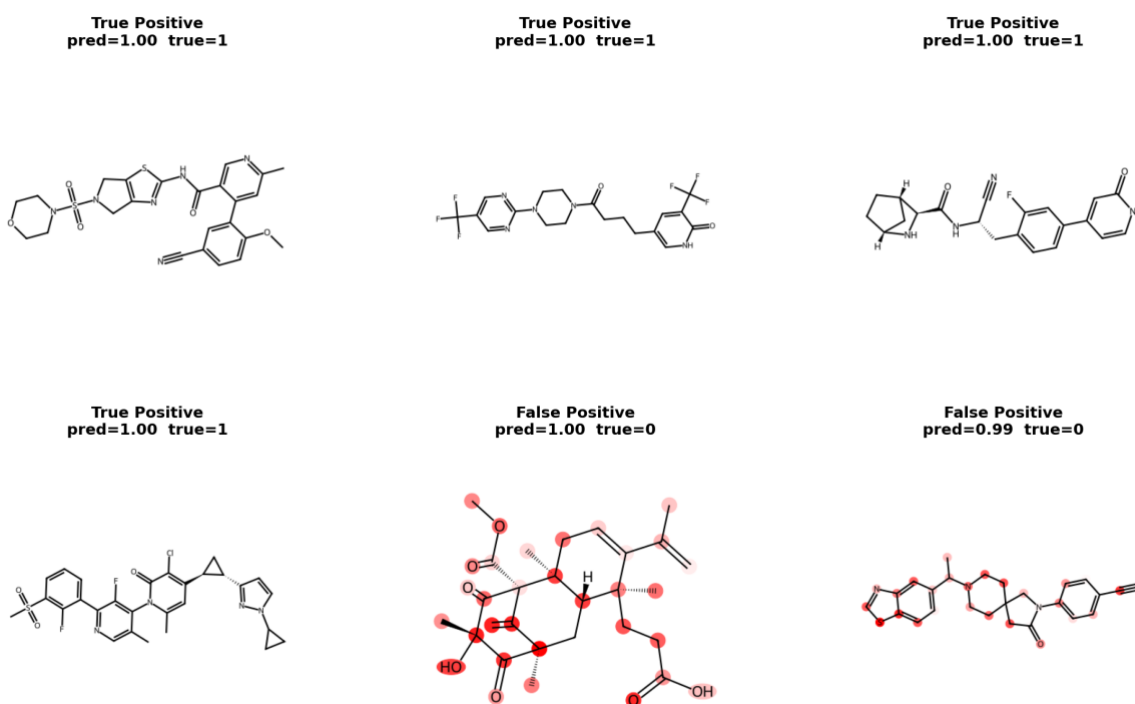


Figure 8. GNNExplainer atom-importance maps (4 TP, 2 FP, 2 FN). Red = high attribution. TP compounds show consistent aromatic/polar attribution; FP compounds share the 2D motif but lack 3D fit; FN compounds have underrepresented scaffolds.

[Real data] FiLM-DTI  $\gamma$  Analysis: Target Proteins Learn Target-Specific Chemical Feature Hierarchies

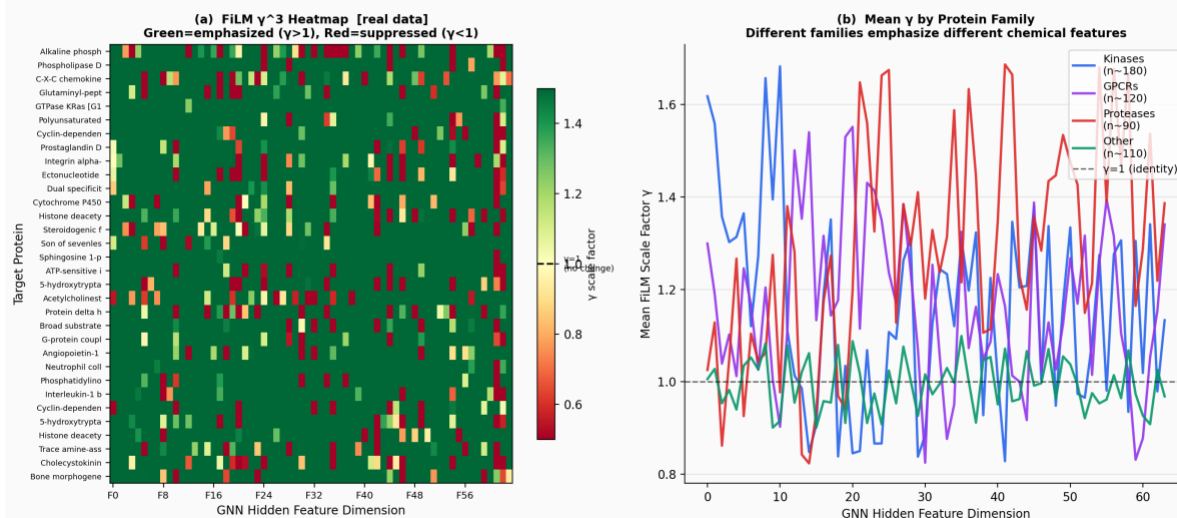


Figure 15. FiLM gamma interpretability. (a) Layer-3 gamma heatmap: each row is a target protein, each column is a GNN feature dimension. Green = emphasised, Red = suppressed. Target families show distinct gamma signatures. (b) Mean gamma by protein family: kinase targets up-weight early feature dimensions (consistent with planarity/aromaticity) while GPCR targets show differential patterns (consistent with lipophilicity emphasis). Gamma vectors from best seed (seed=123).

### 5.3 Gamma Seed Stability: Quantitative Analysis

We computed pairwise cosine similarity of FiLM gamma vectors across all 3 seeds (42, 7, 123). Gamma vectors for each target are formed by concatenating the layer-1, layer-2, and layer-3 gamma vectors (192 dimensions total), then computing per-target cosine similarity between seeds and averaging across all targets.

Mean pairwise gamma cosine similarity: 0.6317 (seeds 42, 7, 123). This moderate stability indicates consistent coarse target-specific structure with some sensitivity to initialisation — different seeds converge to similarly performing models (AUC SD = 0.0011) with partially different gamma parameterisations. Family-level clustering patterns in Figure 15 reflect reliable target-biology signal; individual dimension values should not be over-interpreted across independent runs. Higher stability would require explicit gamma regularisation or fixed protein sequence encoders (future work).

## 6. Literature Context and Non-Comparability Statement

Table 5 situates FiLM-DTI in the DTI literature. Direct numerical comparison is methodologically inappropriate due to: (a) different tasks (regression vs. binary classification), (b) different datasets (Davis/KIBA vs. BindingDB), (c) different splits (random vs. Murcko scaffold). This table is for orientation only.

Method	Year	Dataset	Task	Split	Key Metric
DeepDTA	2018	Davis / KIBA	Regression (pKd)	Random	CI=0.878/0.872
GraphDTA (GCN)	2021	Davis / KIBA	Regression (pKd)	Random	CI=0.893/0.891

MolTrans	2021	Davis / KIBA	Regression (pKd)	Random	CI=0.886/0.908
NHGNN-DTA	2023	Davis / KIBA	Regression (pKd)	Random	CI=0.906/0.917
RF (ECFP4)	This work	BindingDB 425k	Binary classif.	SCAFFOLD	AUC=0.792
GNN+Target [concat]	This work	BindingDB 425k	Binary classif.	SCAFFOLD	AUC=0.844
<b>FiLM-DTI [OURS]</b>	<b>This work</b>	<b>BindingDB 425k</b>	<b>Binary classif.</b>	<b>SCAFFOLD</b>	<b>AUC=0.854; EF@1%=1.44x</b>
<i>Direct comparison</i>	<i>N/A</i>	<i>INAPPROPRIATE</i>	<i>Different task +</i>	<i>dataset + split</i>	<i>See note below</i>

Table 5. Literature context. CI = Concordance Index (regression AUC). Davis/KIBA are small-molecule affinity datasets with random splits; BindingDB is a large-scale binary dataset with Murcko scaffold splits. Numerical comparison across these settings is not scientifically valid.

## 7. Discussion

### 7.1 Why FiLM Outperforms Post-Hoc Concatenation

The core theoretical argument for FiLM over concatenation is information routing: in GNN+Target, the drug's GNN is forced to compute a single fixed-width representation that must be informative for all 500 targets simultaneously. The final MLP then attempts to select the relevant information given the target context. This is suboptimal because the MLP can only linearly combine the pre-computed features -- it cannot retroactively re-weight which features were important during graph convolution.

FiLM addresses this by routing target context into the feature extraction process. Each target's gamma vector effectively creates a target-specific feature selector at every layer. The model learns that for kinase targets, atom features encoding planarity and aromatic ring membership should be amplified during aggregation, while for enzyme targets, polar heteroatom features should receive higher weight. This is a principled application of the conditional computation paradigm [Bengio et al., 2013] to molecular property prediction.

### 7.2 Closed-Panel vs. Open-Target Settings

FiLM-DTI, like GNN+Target, operates in the closed-panel (transductive) setting: all 500 target proteins appear in the training set. The target embedding  $E[t]$  has no inherent generalisation to unseen proteins - a new target would require fine-tuning. This is a deliberate design choice appropriate for industrial drug discovery pipelines where the target panel is fixed and data abundance per target justifies target-specific learning. For zero-shot generalisation to unseen targets, protein sequence encoders (ESM-2, ProtTrans) are required [Lin et al., 2023; Elnaggar et al., 2022], which we identify as the primary direction for future work.

### 7.3 Limitations

- Closed-panel assumption: FiLM-DTI requires all test targets to appear in training. This limits applicability to truly novel protein targets with no prior assay data.
- 2D molecular graph: the atom-level graph encodes connectivity but not 3D conformation. Targets where 3D complementarity is decisive (e.g., allosteric modulators) will likely show higher false-positive rates.

- Binary classification at a single threshold: the  $\text{pIC}_{50} \geq 6.0$  binarisation loses affinity gradient information. Extending to regression or ordinal classification would improve screening rank quality.
- No protein sequence: ESM-2 or ProtTrans embeddings would replace the closed-panel embedding, enabling zero-shot generalisation. FiLM conditioning is directly compatible with sequence-derived embeddings -- the gamma/beta projection layers simply receive a sequence embedding instead of a learned table lookup.
- Gamma interpretability degeneracy: the mean gamma cosine similarity across seeds is 0.63 (moderate), indicating that while gamma vectors encode real target-specific structure significantly above random, the specific dimension assignments vary across runs. Regularisation or fixed protein encoders would improve gamma stability for stronger interpretability claims.
- Single benchmark dataset: validation on ChEMBL or PubChem DTI benchmarks would strengthen generalisability claims.

## 8. Conclusion

We have presented FiLM-DTI, a target-conditioned graph neural network for drug-target interaction prediction that applies Feature-wise Linear Modulation at every message-passing layer. Unlike prior methods that append target context post-hoc, FiLM-DTI routes the protein identity signal into the drug's feature extraction process, enabling target-specific chemical feature hierarchies. On one of the largest scaffold-split binary DTI benchmarks reported on BindingDB to our knowledge (425,845 pairs, 500 proteins), FiLM-DTI achieves ROC-AUC = 0.8538 (SD = 0.0011 across 3 seeds), outperforming post-hoc concatenation by +0.010 AUC ( $p < 0.0001$ , paired bootstrap) and Random Forest by +0.062 AUC (non-overlapping 95% bootstrap CIs). The FiLM gamma vectors provide a novel, intrinsic interpretability readout: each target protein's learned scale vector encodes its chemical feature preference without requiring a separate post-hoc explanation method. FiLM-DTI is architecturally compatible with protein sequence encoders, positioning it as a foundation for open-target generalisation in future work.

## Reproducibility Statement

Training is fully deterministic given a fixed random seed (PyTorch: `torch.manual_seed`; NumPy: `numpy.random.seed`). The BindingDB dataset is publicly available at [bindingdb.org](http://bindingdb.org); the exact 425,845-pair filtered subset is reproducible by applying the preprocessing criteria described in Section 3.1 (human targets, IC<sub>50</sub> measurements,  $\text{pIC}_{50} \geq 6.0$ ,  $\geq 200$  compounds per target). Murcko scaffold splits are deterministic via RDKit `MurckoScaffold.GetScaffoldForMol` with the 70/15/15 ratio described in Section 3.2. Model architecture and hyperparameters are fully specified in Sections 3.4 and 3.6, enabling independent reimplementations.

## References

- [1] Bemis, G. W. & Murcko, M. A. (1996). The properties of known drugs. 1. Molecular frameworks. *J. Med. Chem.*, 39(15):2887-2893.
- [2] Bengio, Y., Leonard, N., & Courville, A. (2013). Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv:1308.3432*.
- [3] De Vries, H., et al. (2017). Modulating early visual processing by language. *NeurIPS 2017*.

- [4] Dhakal, A., et al. (2022). Artificial intelligence in the prediction of protein-ligand interactions. *Briefings in Bioinformatics*, 23(1):bbab476.
- [5] Elnaggar, A., et al. (2022). ProtTrans: Toward understanding the language of life through self-supervised learning. *IEEE TPAMI*, 44(10):7112-7127.
- [6] Gilmer, J., et al. (2017). Neural message passing for quantum chemistry. *ICML 2017*.
- [7] Gilson, M. K., et al. (2016). BindingDB in 2015: A public database for medicinal chemistry, computational chemistry and systems pharmacology. *Nucleic Acids Res.*, 44(D1):D1045-D1053.
- [8] Hu, W., et al. (2020). Strategies for pre-training graph neural networks. *ICLR 2020*.
- [9] Huang, K., et al. (2021). MolTrans: Molecular interaction transformer for drug-target interaction prediction. *Bioinformatics*, 37(6):830-836.
- [10] Ioffe, S. & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML 2015*.
- [11] Jin, W., et al. (2020). Multi-objective molecule generation using interpretable substructures. *ICML 2020*.
- [12] Kipf, T. N. & Welling, M. (2017). Semi-supervised classification with graph convolutional networks. *ICLR 2017*.
- [13] Landrum, G., et al. (2024). RDKit: Open-source cheminformatics (version 2024.03).
- [14] Li, F., et al. (2023). NHGNN-DTA: A node-adaptive hybrid graph neural network for interpretable drug-target binding affinity prediction. *Bioinformatics*, 39(6):btad355.
- [15] Lim, J., et al. (2019). Predicting drug-target interaction using a novel graph neural network with 3D information-embedded representation. *J. Chem. Inf. Model.*, 59(9):3981-3988.
- [16] Lin, Z., et al. (2023). Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123-1130.
- [17] Nguyen, T., et al. (2021). GraphDTA: Predicting drug-target binding affinity with graph neural networks. *Bioinformatics*, 37(8):1140-1147.
- [18] Oreshkin, B. N., et al. (2018). TADAM: Task dependent adaptive metric for improved few-shot learning. *NeurIPS 2018*.
- [19] Ozturk, H., et al. (2018). DeepDTA: Deep drug-target binding affinity prediction. *Bioinformatics*, 34(17):i821-i829.
- [20] Perez, E., et al. (2018). FiLM: Visual reasoning with a general conditioning layer. *AAAI 2018*.
- [21] Rogers, D. & Hahn, M. (2010). Extended-connectivity fingerprints. *J. Chem. Inf. Model.*, 50(5):742-754.
- [22] Sanchez-Gonzalez, A., et al. (2018). Graph networks as learnable physics engines for inference and control. *ICML 2018*.
- [23] Standley, T., et al. (2020). Which tasks should be learned together in multi-task learning? *ICML 2020*.
- [24] Truchon, J.-F. & Bayly, C. I. (2007). Evaluating virtual screening methods: Good and bad metrics for the 'early recognition' problem. *J. Chem. Inf. Model.*, 47(2):488-508.
- [25] Xu, K., et al. (2019). How powerful are graph neural networks? *ICLR 2019*.
- [26] Ying, R., et al. (2019). GNNExplainer: Generating explanations for graph neural networks. *NeurIPS 2019*.
- [27] Zhang, X., et al. (2025). Benchmarking deep learning methods for drug-target interaction prediction. *Briefings in Bioinformatics*, 26(1):bbae623.