INTEGRATING GRAPH NEURAL NETWORKS AND SPECTROPHOTOMETRY FOR pIC50 PREDICTION OF PESTICIDES

<u>Ang</u> Yu Xi Sophie¹, <u>Yap</u> Xiu Huan², <u>Shen</u> Bingquan² ¹Raffles Institution, 1 Raffles Institution Lane, Singapore 575954 ²DSO National Laboratories, 12 Science Park Drive, Singapore 118225

1. Abstract

Accurate pIC50 prediction is vital for assessing the toxicity and potency of chemical compounds, including pesticides. This study leverages computational and experimental approaches to achieve this goal. A Graph Isomorphism Network (GIN) was pretrained on a large PubChem dataset (117,520 molecules) to learn node-level molecular features and fine-tuned on a smaller dataset (14,611 molecules) containing experimentally determined pIC50 values. The model predicts pIC50 values directly from molecular SMILES strings. Validation of the GIN's predictions was conducted using spectrophotometric assays to determine empirical pIC50 values for two pesticides of differing potency. A strong correlation between predicted and experimental values highlights the potential of combining graph-based deep learning for toxicity prediction.

2. Introduction

Molecular property prediction is a key area of computational chemistry, aimed at developing models that map molecular structures to their properties.¹ The accurate prediction of a compound's inhibitory concentration (pIC50) is a cornerstone of toxicological studies, allowing chemical structures, simulation, and physical data to be integrated when predicting health risks and other toxicological information.² The pIC50 value, which quantifies the half maximal inhibitory concentration, is used to determine the potency of a drug against a variety of enzymes or biological targets associated with the pathogenesis of multiple disorders.³ Experimental identification of the bioactivity of each potential peptide is an expensive and time-consuming task.⁴ Deep representation learning has been reported as a promising approach for molecular property prediction, outperforming fixed molecular representations.¹

However, accurately predicting molecular properties poses a significant challenge due to the complex relationships between molecular structures and their properties.⁵ Molecules are typically represented in three ways: fixed representations, including fingerprints and structural keys, that signify the presence of specific structural patterns; linear notations, such as Simplified Molecular Input Line Entry System (SMILES) strings; and molecular graphs.⁶ With the advent of deep learning, various neural networks have been proposed for molecular representation learning, such as convolutional neural networks (CNNs), recurrent neural networks (RNNs) and graph neural networks (GNNs).¹

One major task for AI in drug discovery is molecular property prediction.¹ Graph Isomorphism Networks (GINs) address this challenge by representing molecules as graphs, capturing both local and global features to link chemical topology to biological activity. Furthermore, transfer learning has emerged as an essential technique to address the scarcity of

labelled data and high dimensionality of feature spaces.⁵ By leveraging knowledge gained from related tasks, transfer learning enables models to generalise better to target tasks with limited training data, enhancing prediction accuracy. Thus, to enhance model accuracy, pretraining on large datasets followed by task-specific fine-tuning is crucial.

The EU has established a new chemical assessment paradigm (REACH, enacted in 2007) to evaluate chemicals globally.⁵ In the USA, quantitative structure–activity relationships (QSARs) predictions are used to evaluate two to three thousand chemicals each year and to assess a significant portion of the toxicity information.⁷ However, traditional QSAR models often rely on predefined molecular descriptors (e.g., atom counts, bond types, molecular weight) and linear or nonlinear regression models to establish relationships between molecular features and biological activity.⁸ The coverage and comprehensiveness of metabolic data presents a critical challenge with respect to QSAR modelling⁸, requiring manual recalibration for new datasets or molecular classes.

In this study, we developed a GIN model pretrained on a PubChem toxicity dataset of over 100,000 molecules and fine-tuned it on over 14,000 molecules with known pIC50 values to predict pesticide inhibition. Unlike the static descriptors of QSAR, the GIN directly learns from molecular graph representations and captures topological and relational features more comprehensively. Additionally, the GIN can incorporate new data seamlessly through retraining. Validation through acetylcholinesterase (AChE) spectrophotometric assays revealed strong agreement between predicted and experimental pIC50 values. B

3. Materials and Methods

3.1 Dataset and Data Preprocessing

SMILES Representation: Molecular structures in both datasets were represented using the SMILES notation. To prepare the data for the GIN model, each SMILES string was converted into a graph structure. In this representation, nodes represent the atoms in the molecule and edges represent the bonds between atoms. Bond types (single, double, etc.) were encoded as edge attributes. This transformation enabled the model to process molecular structures in a graph format, capturing both the local connectivity of atoms and the overall topology of the molecule.

Feature extraction: Feature extraction was conducted to provide meaningful inputs to the GIN model at both the node and graph levels. To extract node level features, each atom in the graph was assigned a feature vector encoding its chemical properties, such as atom type, hybridization state and presence of formal charges, to allow the GIN to learn how individual atoms contribute to the molecule's overall activity. In addition, graph-level descriptors such as molecular weight, lipophilicity and topological polar surface area (TPSA) were computed, summarising global molecular properties particularly relevant for pIC50 prediction. Such preprocessing pipeline ensured that the GIN could learn from both the detailed local information at the atom level and the broader structural and physicochemical properties of each molecule.

3.2 Model Architecture and Pretraining

GIN Components: The GIN was chosen for this study due to its ability to distinguish nonisomorphic graphs effectively, making it ideal for molecular graph representation. The key components of the GIN architecture are as follows: (i) the input layer consisting of molecular graphs converted from SMILES strings; (ii) graph convolution layers aggregating information from neighbouring nodes, iteratively updating each node's representation to encode local structural and chemical information; (iii) hidden layers triggering nonlinear transformations using activation functions (e.g. ReLU) and batch normalisation to stabilise training; (iv) the readout layer, a global pooling layer combining node features into a fixed-size graph-level representation; and (v) the task-specific output layer – during pretraining, the output is node-level predictions to learn molecular substructures, while during fine-tuning, the output is a single graph-level target (i.e. pIC50).

Pretraining on the PubChem Dataset: To build a robust feature extractor, the GIN was pretrained on a large toxicity dataset from PubChem via a self-supervised learning task, where the goal was to predict node-level features, such as atomic environments or chemical properties. Using the generative reconstruction technique, we masked the features of a random batch of nodes, forwarded the masked graph through the GIN encoder and reconstructed the masked node features given the node representations of their local sub-graphs⁹. The model parameters were then optimised using gradient descent and a loss function that measured the accuracy of node-level predictions.

Fine-tuning for pIC50 Prediction: After pretraining, the GIN was fine-tuned on a smaller dataset with corresponding pIC50 values. This phase involved task-specific adjustments, where (i) the output layer was modified to predict graph-level instead of node-level properties; (ii) both node-level features (e.g., atomic types) and graph-level descriptors (e.g., molecular weight, lipophilicity, TPSA) were utilised to enhance prediction accuracy; and (iii) the pre-trained weights were used to initialise the GIN, allowing the model to retain its learned molecular features. A regression loss function (e.g., Mean Squared Error) was employed to measure the difference between predicted and actual pIC50 values. The model was fine-tuned for several epochs to reduce training and validation loss. This two-step approach of pretraining and fine-tuning allowed the GIN to combine general molecular feature extraction with task-specific optimisation, significantly improving its accuracy and robustness in predicting pIC50 values.

3.3 Experimental Setup

Colorimetric Determination of Acetylcholinesterase (AChE) Activity: To validate the predictions of the GIN, the toxicities of two pesticides, Profenofos and Dichlorvos, were experimentally determined using an AChE inhibition assay. A preliminary experiment to determine AChE activity was conducted before testing pesticide inhibition to optimise the dilution factor of the AChE enzyme – the goal was to achieve a baseline V_{max} (i.e., maximum reaction rate) of approximately 100 without the addition of pesticides. After attaining the optimal dilution factor of enzyme, the chemical procedure was as follows: (i) preparation: stock solutions of the two pesticides were prepared, serially diluted to create eight concentrations for testing covering a wide range to accurately capture dose-response behaviour; (ii) enzyme incubation: in 96-well plates, triplicates were prepared for each pesticide concentration to ensure reproducibility. Control wells (without pesticide) were included to establish baseline AChE activity. AChE enzyme solution was added to each well containing one of the pesticide concentrations, the mixtures were incubated at 25°C for 15 minutes to allow the pesticides to interact with and inhibit the enzyme; (iii) addition of Ellman's Assay: Ellman's reagent (5,5'-dithiobis-(2-nitrobenzoic acid)) and the substrate acetylcholine were added to the wells. Ellman's assay produces a yellow colour as AChE catalyses the breakdown of acetylcholine, with the intensity proportional to the enzyme's activity¹⁰; (iv) spectrophotometry: a spectrophotometer was used to measure the colour intensity of the solutions and the reaction rates (i.e. V_{max}) at a wavelength of 412 nm, reflecting the enzyme's activity at each pesticide concentration; and (v) repetition for consistency: to enhance the reliability of results, the entire experiment was repeated four times for each pesticide, and the mean values for each measurement were calculated. This comprehensive setup provided a robust dataset for statistical analysis and reduced the influence of experimental variability.

Data Collection: To quantify enzyme inhibition, the following steps were performed: (i) V_{max} Measurement: the spectrophotometer recorded the V_{max} values for each pesticide concentration in triplicate; (ii) averaging V_{max} : the V_{max} values for each concentration were averaged across the triplicates; (iii) enzyme activity calculation: the enzyme activity at each concentration was expressed as a percentage relative to a control (enzyme activity without pesticide) using the formula: Enzyme activity (%) = V_{max} (sample) / V_{max} (control) x 100;

and (iv) percentage inhibition: The percentage inhibition of AChE was calculated as: inhibition (%) = 100 - Enzyme Activity (%). These calculations were repeated for all eight pesticide concentrations for both pesticides.

pIC50 Determination: The pIC50 values of the pesticides were determined using the following approach: (i) data plotting: the percentage inhibition values at each pesticide concentration were entered into PRISM GRAPHPAD software; (ii) sigmoidal curve fitting: a dose-response curve of percentage inhibition against log-transformed pesticide concentration was plotted (*Fig. 1*); (iii) pIC50 determination: the Sigmoidal, 4PL function was used to determine the IC50 value, which was converted to pIC50 using the formula: $pIC50 = -log_{10}(IC50)$. The resulting pIC50 values provided a quantitative measure of each pesticide's inhibitory potency, which was then compared to the GIN's predicted values for model validation.

Separadal, 4FL Xi lagicancentraland Te (7): In a data by the up class, yokak and thinking have data. The space on expression to generative data and the space of the space o		Rules for Initial Values	Default Constraints Tran	nsforms to Report		
Spen - Top - Entrom - (Expected in the Control of Contr	Signu Tip: H -This e X: log Y: Rec Top a logIC5 HillSlo	aidal, 4PL, X is leg(c. X is not already the log or quastion is equivalent to: al dose or concentration prome, discharging the discharging of discharging the discharging of the series log undit as X as Stope factor or Hill sto	nncentration) f dose, go back and transf logidose) vs. response (va increases me units ar Y. spe, unifies.	sm your data. iiable slope)		
	ipan = r'=Botts	Top · Bottom m + (Top-Bottom)/(1+10	'((LogiCSOX((HillSlope))	^	Tip-	_

Fig. 1. Built-in equation for Sigmoidal, 4Pl function of the PRISM GRAPHPAD software

4. Results

4.1 GIN Prediction Results

Prediction accuracy:



Fig. 2. Scatter plots of predicted vs actual pIC50 values using simpler models: (from left to right) simple linear regression, random forest regression and GIN

Model	Mean squared error
Simple linear regression	1.4881
Random forest regression	0.5789
Graph isomorphism network	0.5698

Fig. 3. Table of	^r values of mean	squared error for	r each mode	l employed
------------------	-----------------------------	-------------------	-------------	------------

Epoch 1, Training Loss: 0.7303854137496218
Epoch 1, Validation Loss: 0.6841405176598093
Model saved at epoch 1 with validation loss: 0.6841405176598093
Epoch 2, Training Loss: 0.580572730042244
Epoch 2, Validation Loss: 0.7463098336821017
Epoch 3, Training Loss: 0.523014206052478
Epoch 3, Validation Loss: 0.6647760103578153
Model saved at epoch 3 with validation loss: 0.6647760103578153
Epoch 4, Training Loss: 0.4844845746710001
Epoch 4, Validation Loss: 0.6399535050858622
Model saved at epoch 4 with validation loss: 0.6399535050858622
Epoch 5, Training Loss: 0.4589613482274644
Epoch 5, Validation Loss: 0.7705744919569596
Epoch 6, Training Loss: 0.4385773071337267
Epoch 6, Validation Loss: 0.7156230910964634
Epoch 7, Training Loss: 0.42326210919625123
Epoch 7, Validation Loss: 0.6127975868142169
Model saved at epoch 7 with validation loss: 0.6127975868142169
Epoch 8, Training Loss: 0.4083780018521137
Epoch 8, Validation Loss: 0.6275735769582831
Epoch 9, Training Loss: 0.3965139038901511
Epoch 9, Validation Loss: 0.6954140481741532
Epoch 10, Training Loss: 0.38277937610292695
Epoch 10, Validation Loss: 0.5683993489845939
Model saved at epoch 10 with validation loss: 0.5683993489845939

Fig. 4. Training and validation losses over 10 training epochs, with the GIN model saved at the best validation loss

new_molecule1_smiles = "COP(=0)(0C)0C=C(C1)C1"
predicted_pIC50 = predict_pIC50(new_molecule1_smiles, model, data_mean, data_std)
print(f"Predicted pIC50 for {new_molecule1_smiles}: {predicted_pIC50}")

new_molecule2_smiles = "CCCSP(=0)(OCC)OC1=C(C=C(C=C1)Br)C1"
predicted_pIC50 = predict_pIC50(new_molecule2_smiles, model, data_mean, data_std)
print(f"Predicted pIC50 for {new_molecule2_smiles}: {predicted_pIC50}")

Predicted pIC50 for COP(=0)(0C)0C=C(C1)C1: 3.00315197293448 Predicted pIC50 for CCCSP(=0)(0CC)0C1=C(C=C1)Br)C1: 4.224558182399754 Fig. 5. GIN's predicted pIC50 values, where molecule 1 represents Dichlorvos and molecule 2 represents

Profenofos

4.2 Experimental Results

Inhibition data:

famala	Cons (nM)	Los cons	MAR	Manan	Maan Verav	Mean activity (a.v.)	Inhibition (M)	Sample	Conc (nM)	Log conc	Well	Vmax	Mean Vmax	Mean activity (a.u.)	Inhibition (%)
Sample	conc (mvi)	Log conc	wei	VIIIAA	mean vinax	wear activity (a.u.)	innoidion (Ay		0	-3	A1	73.99	70.67	100	0
1	U	-3	A1	71.82	73.853	100	0	1			A2	68.16			
			A2	75.24							A3	69.86			
	100		A3	74.5	70.110	07.01000010	0.050004570		100	2	B1	72.77	69.38	98.1746144	1.825385595
2	100	2	81	70.45	72.113	97.64396842	2.356031576	2			82	67.37			
2 2 3			82	72.08							B3	68			
	250		63	/3.81		00 10001700	0.01100007		250	2.4	C1	61.86	59.173	83.73142776	16.26857224
3	250	2.4	0	65.11	66.607	90.18861793	9.81138207	3			C2	58.57			
				00.12							C3	57.09			
	500		6	66.59	67.040	77.1007000	00 5040445		500	2.7	D1	48.45	46.81	66.23744163	33,76255837
4	500	2.7	DI	55.51	57.213	//.468/555	22.5312445	4			D2	45.95			
4			02	58.32							D3	46.03			
	4000		03	57.81	10 707	00.00050.007	00.00747500		1000	3	E1	36.67	35 777	50 6254422	49 3745578
	1000	3	El	47.42	48.767	66.03252407	33.96747593	5	1000		F2	34.38	55.111	00.0204422	40.0140010
5			E2	49.55							F3	36.28			
			E3	49.33					6691	2.92	E1	1.64	2 202	4 650695964	05 24021414
	6691	3.83	F1	9.21	6.843	9.265703492	90.73429651	6	0091	3.65	62	4.45	3.233	4.005000004	50.34031414
6			F2	6.7				0			F2	4.40			
			R	4.62					19993	4.12	F3	3.79	1.90	3 631051333	07 26904969
-	13382	4.13	61	1.097	1.196	1.619433198	98.3805668		13302	4.15	61	3.30	1.00	2.031931323	97.30004000
7			G2	1.09				/			62	1.04			
			G3	1.4							63	1.16			
	26764	4.43	H1	-0.09	0.573	0.775865571	99.22413443	8	26764	4.43	HI	3.03	0.51	0.721664072	99.27833593
8			H2	0.22							HZ	-0.6			
			H3	1.59							H3	-0.9			
Sample	Conc (nM)	Log conc	Well	Vmax	Mean Vmax	Mean activity (a.u.)	Inhibition (%)	Sample	Conc (nM)	Log conc	Well	Vmax	Mean Vmax	Mean activity (a.	a.) Inhibition (9
	0	-3	A1	108.21	102.15	100	0		0	-3	A1	78.53	76.283	100	0
1			A2	102.1				1			A2	76.01			
			A3	96.14							A3	74.31			
	100	2	B1	99.05	100.023	97,91776799	2.082232012		100	2	B1	69.68	70.907	92.95255824	7.04744176
2			B2	101.73				2			B2	72.58			
			B3	99.29							B3	70.46			
	250	2.4	C1	86.63	87.463	85.62212433	14.37787567		250	2.4	C1	60.41	63.92	83.79324358	16.2067564
3			C2	88.67				3			C2	68.74			
			C3	87.09							C3	62.61			
	500	2.7	D1	76.64	78.45	76,79882526	23.20117474		500	2.7	D1	49.88	54.26	71,12987166	28.8701283
4			D2	80.88				4			D2	58.97			
			D3	77.83							D3	53.93			
	1000	3	E1	68.85	64.137	62,78707783	37.21292217		1000	3	E1	34.18	40.537	53,14028027	46.8597197
5		-	F2	63.52				5			E2	42.96			
			F3	60.04				-			E3	44.47			
	6691	3.83	F1	17 71	14 347	14 04503182	85 95496818		26764	4.43	F1	3.42	4 514	5 91743901	94 0825605
6	0051	5.05	E2	15.14	14.047	14.04000102	00.00400010	6			F2	5 843			
5			B	10.19				1			F3	4 279			
	13382	4.13	61	8.214	6.666	6.525697504	93 4743025		33456	4.52	G1	-0.7	2.647	3.469973651	96,5300267
7	LUJUE		62	7.07	0.000	0.020007.004	30.4140020	7			62	5.66	2.041	0.100010001	00.0000200
-				7.07			1	-			01	0.00			
			63	4 714							63	2.98			
	26764	4.43	G3 H1	4.714	4 502	4 407244249	05 50275575		40147	4.6	G3	2.98	2 987	3 015682303	96.0843176
	26764	4.43	G3 H1	4.714 3.614	4.502	4.407244249	95.59275575		40147	4.6	G3 H1	2.98	2.987	3.915682393	96.0843176
8	26764	4.43	G3 H1 H2	4.714 3.614 4.491	4.502	4.407244249	95.59275575	8	40147	4.6	G3 H1 H2	2.98 -0.96 5.79	2.987	3.915682393	96.0843176

Fig. 6. Experimental data for **Profenofos** across four experiments, read from left to right, top to bottom, with log concentration (x-axis) and percentage inhibition (y-axis) boxed

								Sample	Conc (nM)	Log conc	Well	Vmax	Mean Vmax	Mean activity (a.u.	.) Inhibition (%)
Sample	Conc (nM)	Log conc	Well	Vmax	Mean Vmax	Mean activity (a.u.)	Inhibition (%)		0	-3	A1	66.41	65.153	100	0
	0	-3	Al	72.07	68.817	100	0	1			A2	66.27			
1			A2	68.46							A3	62.78			
			A3	65.92					100	2	B1	59.33	56,957	87.42037972	12.57962028
2	199	2	81	68.06	65.003	94.4577648	5.542235204	2			B2	55.24			
			BZ	64.15							B3	56.3			
			83	62.8					250	2.4	C1	46.22	44 647	68 52639172	31.47360828
	497.5	2.4	61	62.22	60.343	87.68618219	12.31381781	3			02	43.96			
3			62	61.19				-			C3	43.76			
			C3	57.62					500	2.7	D1	32.23	30 757	47 20734272	52 79265728
4	995	2.7	D1	50.95	51.05	74.18225148	25.81774852	4			D2	30.81	00.101	Theorement	OL. I OLOOT LO
4			D2	53.59							03	20.23			
			03	48.61					1000	3	E1	16.7	14.22	21 07002207	78.02006603
-	1990	3	E1	39.38	37.927	55.11283549	44.88716451	6	1000	-	F2	12.10	14.32	21.07803387	78.02050003
5			EZ	37.69				-			62	13.13			
			E3	36.71					2000	2.2	63	13.07	0.70	4.040500505	00 70740040
6	112256.68	3.75	F1	-0.6	0.18	0.261563277	99.73843672		2000	3.3	11	1.2	0.79	1.212530505	96.76740949
			F2	NoFit				6			12	1.471			
			F3	0.96							F3	-0.3			
	16885.03	3.93	G1	1.131	1.02	1.4821919	98.5178081	7	2500	3.4	G1	0	0	0	100
7			G2	1.02							62	0			
			G3	0.91							G3	0			
	22513.37	4.05	Н1	NoFit	NoFit				4000	3.6	H1	0	1.8	2.762727733	97.23727227
8			H2	NoFit				8			H2	1.8			
			H3	NoFit							H3	0			
Sample	Conc (nM)	Log conc	Well	Vmax	Mean Vmax	Mean activity (a.u.	Inhibition (%)	Sample	Conc (nM)	Log conc	Well	Vmax	Mean Vmax	Mean activity (a.u.)	Inhibition (%)
	0	-3	A1	55.59	54.77	100	0	Jumple	0	-2	41	102.564	102 962	100	0
1			Δ2	54.81					0	~	~	102.004	103.002	100	0
								1			42	103 264			
			A3	53.91				1			A2	103.264			
	100	2	A3 B1	53.91 50.67	49.437	90.26291766	9.737082344	1	100	2	A2 A3 B1	103.264 105.757 97.764	98 119	94 47054746	5 529452543
2	100	2	A3 B1 B2	53.91 50.67 49.24	49.437	90.26291766	9.737082344	1	100	2	A2 A3 B1 B2	103.264 105.757 97.764 95.914	98.119	94.47054746	5.529452543
2	100	2	A3 B1 B2 B3	53.91 50.67 49.24 48.4	49.437	90.26291766	9.737082344	2	100	2	A2 A3 B1 B2 B3	103.264 105.757 97.764 95.914 100.679	98.119	94.47054746	5.529452543
2	250	2	A3 B1 B2 B3 C1	53.91 50.67 49.24 48.4 38.77	49.437 38.397	90.26291766	9.737082344 29.89410261	2	100	2	A2 A3 B1 B2 B3 C1	103.264 105.757 97.764 95.914 100.679 83.086	98.119	94.47054746	5.529452543
2	250	2	A3 B1 B2 B3 C1 C2	53.91 50.67 49.24 48.4 38.77 38.44	49.437 38.397	90.26291766	9.737082344 29.89410261	2	100	2	A2 A3 B1 B2 B3 C1 C2	103.264 105.757 97.764 95.914 100.679 83.086 79.85	98.119 86.974	94.47054746 83.73996264	5.529452543
2 3	100	2	A3 B1 B2 B3 C1 C2 C3	53.91 50.67 49.24 48.4 38.77 38.44 37.98	49.437 38.397	90.26291766 70.10589739	9.737082344 29.89410261	2	250	2	A2 A3 B1 B2 B3 C1 C2 C3	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986	98.119 86.974	94.47054746 83.73996264	5.529452543 16.26003736
2 3	250	2 2.4 2.7	A3 B1 B2 B3 C1 C2 C3 D1	53.91 50.67 49.24 48.4 38.77 38.44 37.98 28.64	49.437 38.397 27	90.26291766 70.10589739 49.29706043	9.737082344 29.89410261 50.70293957	2	250	2	A2 A3 B1 B2 B3 C1 C2 C3 D1	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986 71.771	98.119 86.974 70.838	94.47054746 83.73996264	5.529452543 16.26003736
2 3 4	250	2 2.4 2.7	A3 B1 B2 B3 C1 C2 C3 D1 D2	53.91 50.67 49.24 48.4 38.77 38.44 37.98 28.64 27.33	49.437 38.397 27	90.26291766 70.10589739 49.29706043	9.737082344 29.89410261 50.70293957	2	100 250 500	2.4	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764	98.119 86.974 70.838	94.47054746 83.73996264 68.20396295	5.529452543 16.26003736 31.79603705
2 3 4	100 250 500	2 2.4 2.7	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3	53.91 50.67 49.24 48.4 38.77 38.44 37.98 28.64 27.33 25.03	49.437 38.397 27	90.26291766 70.10589739 49.29706043	9.737082344 29.89410261 50.70293957	1 2 3 4	100 250 500	2 2.4 2.7	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D2	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764 71.070	98.119 86.974 70.838	94.47054746 83.73996264 68.20396295	5.529452543 16.26003736 31.79603705
2 3 4	100 250 500	2 2.4 2.7 3	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1	53.91 50.67 49.24 48.4 38.77 38.44 37.98 28.64 27.33 25.03 13.15	49.437 38.397 27 13.29	90.26291766 70.10589739 49.29706043 24.26510864	9.737082344 29.89410261 50.70293957 75.73489136	1 2 3 4	100 250 500	2	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764 71.979 51.307	98.119 86.974 70.838 49.979	94.47054746 83.73996264 68.20396295 48.12058308	5.529452543 16.26003736 31.79603705 51.87941692
2 3 4 5	100 250 500 1000	2 2.4 2.7 3	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2	53.91 50.67 49.24 48.4 38.77 38.44 37.98 28.64 27.33 25.03 13.15 14.54	49.437 38.397 27 13.29	90.26291766 70.10589739 49.29706043 24.26510864	9.737082344 29.89410261 50.70293957 75.73489136	1 2 3 4 5	100 250 500	2 2.4 2.7 3	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 F2	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764 71.979 51.307 47.729	98.119 86.974 70.838 49.979	94.47054746 83.73996264 68.20396295 48.12058308	5.529452543 16.26003736 31.79603705 51.87941692
2 3 4 5	100 250 500	2 2.4 2.7 3	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3	53.91 50.67 49.24 48.4 38.77 38.44 37.98 28.64 27.33 25.03 13.15 14.54 12.18	49.437 38.397 27 13.29	90.26291766 70.10589739 49.29706043 24.26510864	9.737082344 29.89410261 50.70293957 75.73489136	1 2 3 4 5	100 250 500	2 2.4 2.7 3	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764 71.979 51.307 47.729 50.9	98.119 86.974 70.838 49.979	94.47054746 83.73996264 68.20396295 48.12058308	5.529452543 16.26003736 31.79603705 51.87941692
2 3 4 5	100 250 500 1000	2 24 27 3	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1	53.91 50.67 49.24 48.4 38.877 38.44 37.99 28.64 27.33 25.03 13.15 14.54 12.18 4.92	49.437 38.397 27 13.29 5.62	90.26291766 70.10589739 49.29706043 24.26510864 10.26109184	9.737082344 29.89410261 50.70293957 75.73489136 89.73890816	1 2 3 4 5	100 250 500 1000	2 2.4 2.7 3 3.75	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764 71.979 51.307 47.729 50.9 8.871	98.119 86.974 70.838 49.979 8.119	94,47054746 83,73996264 68,20396295 48,12058308 7,817103464	5.529452543 16.26003736 31.79603705 51.87941692 92.18289654
2 3 4 5 6	100 250 500 1000 2000	2 2.4 2.7 3 3.3	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2	53.91 50.67 49.24 48.4 38.77 38.44 37.98 28.64 27.33 25.03 13.15 14.55 14.55 12.18 4.92 6.21	49.437 38.397 27 13.29 5.62	90.26291766 70.10589739 49.29706043 24.26510864 10.26109184	9.737082344 29.89410261 50.70293957 75.73489136 89.73890816	1 2 3 4 5	100 250 500 1000 5657	2 2.4 2.7 3 3.75	A2 A3 81 82 83 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 E2	103.264 105.767 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764 71.979 51.307 47.729 55.9 8.871 7.586	98.119 86.974 70.838 49.979 8.119	94.47054746 83.73996264 68.20396295 48.12058308 7.817103464	5.529452543 16.26003736 31.79603705 51.87941692 92.18289654
2 3 4 5 6	100 250 500 1000 2000	2 24 27 3 3.3	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2 F3	53.91 50.67 49.24 48.4 38.877 38.44 27.33 25.03 13.15 14.54 12.18 4.92 6.21 5.73	49.437 38.397 27 13.29 5.62	90.26291766 70.10589739 49.29706043 24.26510864 10.26109184	9.737082344 29.89410261 50.70293957 75.73489136 89.73890816	1 2 3 4 5 6	100 250 500 1000 5657	2 2.4 2.7 3 3.75	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F1 F2 F3	103.264 105.757 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764 71.979 68.764 71.979 50.9 8.871 7.586 7.9	98.119 86.974 70.838 49.979 8.119	94.47054746 83.73996264 68.20396295 48.12058308 7.817103464	5.529452543 18.26003736 31.79603705 51.87941692 92.18289654
2 3 4 5 6	100 250 500 2000 2000	2 2.4 2.7 3 3.3 3.3	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2 F3 G1	53.91 50.67 49.24 48.4 38.77 38.44 37.98 28.64 27.33 25.03 13.15 14.54 12.18 4.92 6.21 5.73 -0.03	49.437 38.397 27 13.29 5.62 3.438	90.26291766 70.10589739 49.29706043 24.26510864 10.26109184 6.277159029	9.737082344 29.89410261 50.70293957 75.73489136 89.73890816 93.72284097	1 2 3 4 5 6	100 250 500 1000 5657 8485	2 2.4 2.7 3 3.75	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2 F3 G1	103.264 105.767 97.764 95.914 100.679 83.086 79.85 97.986 71.771 68.764 71.979 51.307 47.729 50.9 8.871 7.586 7.9 8.864	98.119 86.974 70.838 49.979 8.119 4.698	94.47054746 83.73996264 68.20396295 48.12058308 7.817103464 4.623309776	5.529452543 16.26003736 31.79603705 51.87941692 92.18289654 95.47650022
2 3 4 5 6 7	100 250 500 2000 2500	2 24 2.7 3 3.3 3.4	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2 F3 G1 G2	53.91 50.67 49.24 48.4 38.87 28.64 27.33 25.03 13.15 14.54 4.92 6.21 5.73 -0.03 6.08	49.437 38.397 27 13.29 5.62 3.438	90.26291766 70.10589739 49.29706043 24.26510864 10.26109184 6.277159029	9.737082344 29.89410261 50.70293957 75.73489136 89.73890816 93.72284097	1 2 3 4 5 6	100 250 500 1000 5657 8485	2 2.4 2.7 3 3.75 3.93	A2 A3 B1 B2 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2 F3 G1 G2	103.264 105.767 97.764 95.914 100.679 83.086 71.871 83.086 71.871 68.764 71.979 51.307 47.729 50.9 8.871 7.586 7.9 8.864 8.864	98.119 86.974 70.838 49.979 8.119 4.698	94.47054746 83.73996284 68.20396295 48.12058308 7.817103464 4.523309776	5.529452543 16.26003738 31.79603705 51.87941692 92.18289654 95.47669022
2 3 4 5 6 7	100 250 500 2000 2500	2 2.4 2.7 3 3.3 3.3 3.4	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2 F3 G1 G3	53.91 50.67 49.24 48.4 38.77 38.84 27.33 25.03 13.15 14.54 12.18 4.92 6.21 5.73 6.03 6.08 6.08	49.437 38.397 27 13.29 5.62 3.438	90.26291766 70.10569739 49.29706043 24.26510864 10.26109184 6.277159029	9.737082344 29.89410261 50.70293957 75.73489136 89.73890816 93.72284097	1 2 3 4 5 6 7	100 250 500 1000 5657 8485	2 2.4 2.7 3 3.75 3.93	A2 A3 B1 B2 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2 F3 G1 G2 G3	103.264 105.767 97.764 95.914 100.679 83.086 79.985 97.986 71.771 68.764 71.979 51.307 47.729 55.9 8.8671 7.586 7.9 8.864 -0.786 6.414	98.119 86.974 70.838 49.979 8.119 4.698	94.47054746 83.73996264 68.20396295 48.12058308 7.817103464 4.523309776	5.529452543 16.26003736 31.79603705 51.87941692 92.18289654 95.47669022
2 3 4 5 6 7	100 250 500 2000 2500 4000	2 2.4 2.7 3 3.3 3.4 3.6	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 F1 F2 F3 G1 G2 G3 H1	53.91 50.67 49.24 48.4 38.87 28.84 27.33 25.03 13.15 14.54 12.18 4.92 6.21 5.73 -0.03 6.08 4.264 NoFit	49.437 38.397 27 13.29 5.62 3.438	90.26291766 70.10589739 49.29706043 24.26510864 10.26109184 6.277159029 2.702209238	9.737082344 29.89410261 50.70293957 75.73489136 89.73890816 93.72284097 97.29779076	1 2 3 4 5 6 7	100 250 500 1000 5657 8485	2 2.4 2.7 3 3.75 3.93 4.05	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D2 D3 E1 E2 E3 F1 F2 F3 G1 G2 G3 H1	103.264 105.767 97.764 95.914 100.679 83.086 71.978 68.764 71.979 68.764 71.979 51.307 47.729 55.9 8.871 7.586 7.9 8.8464 -0.786 6.414 9.921	98.119 86.974 70.838 49.979 8.119 4.698	94.47054746 83.73996284 68.20396285 48.12058308 7.817103464 4.523309776 2.038514667	5.529452543 18.26003736 31.79603705 51.87941692 92.18289654 95.47669022 97.06148543
2 3 4 5 6 7	100 250 500 2000 2000 4000	2 2.4 2.7 3 3.3 3.4 3.6	A3 B1 B2 B3 C1 C2 C3 D1 D2 D3 E1 E2 E3 F1 F2 F3 G1 G2 G3 H1 H2	53.91 50.67 49.24 48.4 38.77 38.84 27.33 25.03 13.15 14.54 12.18 4.92 6.21 5.73 6.03 6.03 6.08 4.264 NoFit 2.95	49.437 38.397 27 13.29 5.62 3.438 1.48	90.26291766 70.10589739 49.29706043 24.26510864 10.26109184 6.277159029 2.702209239	9.737082344 29.89410261 50.70293957 75.73489136 89.73890816 93.72284097 97.29779076	1 2 3 4 5 6 7	100 250 500 1000 5657 8485 11313	2 2.4 2.7 3 3.75 3.93 4.05	A2 A3 B1 B2 B3 C1 C2 C3 D1 D2 D2 D3 E1 E2 E3 F1 F2 F3 G1 G2 G3 H1	103.264 105.767 97.764 95.914 100.679 83.086 79.985 97.986 71.771 68.764 71.979 51.307 47.729 50.9 8.871 7.586 7.9 8.864 -0.786 6.414 3.021	98.119 86.974 70.838 49.979 8.119 4.698 3.052	94.47054746 83.73996264 68.20396295 48.12058308 7.817103464 4.523309776 2.938514567	5.529452543 16.26003736 31.79603705 51.87941692 92.18289654 95.47669022 97.06148543

Fig. 7. Experimental data for Dichlorvos across four experiments, read from left to right, top to bottom, with log concentration (x-axis) and percentage inhibition (y-axis) boxed

Sigmoidal fit:



Fig. 8. Sigmoidal curves of percentage AChE inhibition against log-transformed concentration of **Profenofos** across four experiments, as plotted on PRISM

Best-fit values	Expt 1	Expt 2	Expt 3	Expt 4	Mean					
Тор	104.5	102.1	105.4	99.58	102.895					
Bottom	0.3186	-1.396	-0.7566	-0.05916	-0.47329					
LogIC50	3.215	2.978	3.227	3.055	3.11875					
HillSlope	1.254	1.224	1.036	1.072	1.1465					
IC50	1642	951.2	1686	1134	1353.3					
Span	104.2	103.5	106.2	99.64	103.385					
	Goodness of Fit									
Degrees of Freedom	4	4	4	4						
R square	0.9984	0.9977	0.9988	0.999	0.998475					
Absolute Sum of Squares	22.05	29.52	14.98	12.34						
Sy.x	2.348	2.716	1.935	1.757						

Fig. 9. Statistical analysis of critical values, particularly pIC50 (in red) and R² (in blue), for Profenofos



Fig. 10. Sigmoidal curves of percentage AChE inhibition against log-transformed concentration of **Dichlorvos** across four experiments, as plotted on PRISM

Best-fit values	Expt 1	Expt 2	Expt 3	Expt 4	Mean						
Тор	101.7	101.9	103.9	104.2	102.925						
Bottom	1.092	0.686	-0.1098	-0.06993	0.3995675						
LogIC50	3.07	2.96	3	2.996	3.0065						
HillSlope	1.356	1.397	1.37	1.233	1.339						
IC50	1176	911.6	1001	989.7	1019.575						
Span	100.6	101.2	104	104.3	102.525						
	Goodness of Fit										
Degrees of Freedom	4	4	4	4							
R square	0.999	0.9993	0.9997	0.9995	0.9994						
Absolute Sum of Squares	13.28	7.926	3.662	6.321							
Sy.x	1.822	1.408	0.9568	1.257							

Fig. 11. Statistical analysis of critical values, particularly pIC50 (in red) and R² (in blue), for Dichlorvos

4.3 Comparison of GIN Predictions and Experimental pIC50 Values



Fig. 12. Bar graph of predicted vs experimental pIC50 values for Profenofos and Dichlorvos

5. Discussion

5.1 Interpretation of Model Performance

Strengths: The GIN demonstrated relatively strong performance in predicting pIC50 values from molecular SMILES strings, highlighting the model's ability to capture the relationship between molecular structure and inhibitory potency. The pretraining on a large toxicity dataset ensured the model learned generalisable molecular features, while fine-tuning on the smaller dataset of labelled pIC50 values enabled it to specialise in this specific task.

Challenges: Despite the model's success, certain challenges were identified: (i) potential overfitting: despite cross validation, the smaller fine-tuning dataset may limit generalisation and increase the risk of overfitting, especially if the model becomes too specialised to the training data; (ii) prediction discrepancies: variability in enzyme activity or spectrophotometric readings, compounded by the training dataset's limited chemical diversity, may contribute to differences between predicted and experimental pIC50 values; (iii) rare or unique molecular structures: underrepresented or absent structural motifs in the training data may reduce prediction accuracy for complex molecules; (iv) simplified assumptions in computational modelling: while GINs excel at capturing molecular topology and physicochemical properties, certain biochemical factors influencing pIC50 values, such as solvent effects, metabolic pathways or protein-ligand dynamics, are not explicitly modelled, contributing to prediction errors; and (v) substructure representation: GINs may struggle with complex substructures, such as large aromatic systems with multiple halogen and functional group substitutions, which present intricate patterns of interactions that can be difficult to fully capture.

5.2 Experimental verification

Agreement between predicted and experimental values: The comparison between the GIN model's predicted pIC50 values and those experimentally-obtained through enzyme inhibition assays provides insights into the model's performance and its ability to generalize from the training data to unseen test compounds. In this section, we analyse the correlation between the predicted and experimental pIC50 values for Dichlorvos and Profenofos — which were used to validate the model. For Dichlorvos, the GIN model demonstrated a high degree of accuracy in predicting the pIC50 value, with a difference of 0.00330 (*Fig. 12*). The empirically determined pIC50 was found to be closely aligned with the model's predicted value, indicating that the model was able to effectively capture the relationship between the SMILES representation of Dichlorvos and its toxicity. In contrast, for Profenofos, the discrepancy between the predicted and experimental

values for Profenofos was more pronounced, with a larger difference of 1.10581 (*Fig. 12*), which suggests that the model struggled with the complex molecular features of this compound.

Evaluation: Dichlorvos achieved an R^2 value of 0.9994, indicating an almost perfect fit between the experimental data and the sigmoidal curve. This high R^2 suggests that the enzyme inhibition data followed a consistent pattern with minimal variability, making it easier for the GIN model to predict the pIC50 value accurately. Profenofos, in contrast, achieved an R^2 of 0.9985, still high but slightly lower than Dichlorvos. This minor drop in R^2 reflects a slight increase in variability in the experimental data.

A comparison between the predicted and experimental pIC50 values revealed a stronger Dichlorvos (COP(=O)(OC)OC=C(Cl)Cl) than Profenofos correlation (CCCSP(=O)(OCC)OC1=C(C=C(C=C1)Br)Cl). The observed discrepancy in model accuracy may stem from several factors: (i) chemical complexity: Dichlorvos is a simpler molecule with a more straightforward structure, consisting of a phosphate group, two ester linkages and a conjugated alkene. Meanwhile, Profenofos' bulkier, more intricate structure, with a large aromatic ring system (i.e. benzene with bromine and chlorine substituents) attached to a phosphorothioate group adds more variability to the molecular structure, including additional aromaticity and halogenation, which may challenge the model's ability to capture binding dynamics; (ii) training data bias: the pretraining dataset may have overrepresented simpler molecules like Dichlorvos and underrepresented complex aromatic systems, reducing accuracy for molecules like Profenofos; (iii) interaction dynamics: while Dichlorvos, being a relatively smaller and more linear molecule, might have more predictable interactions, Profenofos' steric hindrance and complex binding dynamics may exceed the model's current capacity to account for factors like solubility and metabolic transformations that exert a greater influence on the biological activity of Profenofos.

5.3 Implications and Future Work

Room for improvement of GIN model: The lower accuracy for Profenofos is likely due to its complex structure, underrepresented in the training data and the GIN's difficulty handling aromatic and halogenated features. Future improvements could include expanding the dataset with more complex molecules, enhancing feature extraction for subtle structural details, refining the model for nuanced interactions, and optimising training methods to improve accuracy and generalisability.

Potential real-world applications: The GIN model has broad real-world potential in drug discovery, toxicology and environmental science. It can predict the efficacy and toxicity of drug candidates, assess pesticide toxicity to protect non-target organisms and aid in hazard classification and risk assessment for industrial chemicals, pharmaceuticals, and pollutants.

6. Conclusion

Main findings: This study highlights the GIN's ability to predict pIC50 values from SMILES, bridging computational modeling and experimental validation in toxicity assessment. AChE inhibition assays confirmed strong prediction accuracy for Dichlorvos, with more variability for Profenofos. While reliable for simpler molecules, performance declined with increasing complexity of chemical structures, emphasising the need for diverse data and better feature representation. The research demonstrates GINs' potential in drug discovery, toxicology and environmental science.

Acknowledgements

I would like to express my deepest gratitude to Bingquan Shen from the Defence Science Organisation (DSO) for his invaluable guidance, insightful feedback and for making this project possible. I would also like to thank Yap Xiu Huan from DSO for her kind supervision and support through this research. Special thanks to Shen Xiaoting at the Chemical, Toxins and Radionuclear (CTRN) programme for her assistance with the spectrophotometric assays, experimentation and data collection. Lastly, I am grateful to my family and friends for their unwavering support throughout this journey. This research was supported by DSO.

References

- Deng, J., Yang, Z., Wang, H. et al. A systematic study of key elements underlying molecular property prediction. Nat Commun 14, 6395 (2023). https://doi.org/10.1038/s41467-023-41948-6
- 2. Luechtefeld T, Hartung T. Computational approaches to chemical hazard assessment. Altex. 2017;34:459–478. doi: 10.14573/altex.1710141
- 3. Thakur, A. (2022). PIC50: an open source tool for interconversion of PIC50 values and IC50 for efficient data representation and analysis. ResearchGate. https://doi.org/10.1101/2022.10.15.512366
- Singh, D., Mahadik, A., Surana, S., & Arora, P. (2022). Proteochemometric method for PIC50 prediction of flaviviridae. BioMed Research International, 2022, 1–7. <u>https://doi.org/10.1155/2022/7901791</u>
- Jie Shen and Christos A. Nicolaou. Molecular property prediction: recent trends in the era of artificial intelligence. Drug Discovery Today: Technologies, 32-33:29–36, 2019. ArtificialIntelligence
- 6. David, L., Thakkar, A., Mercado, R. & Engkvist, O. Molecular representations in ai-driven drug discovery: a review and practical guide. J. Cheminformatics 12, 1–22 (2020)
- 7. Ministry of Environment, Republic of Korea (2001) Study on improvement of new chemical substance hazard assessment system
- Cherkasov A, Muratov EN, Fourches D, Varnek A, Baskin II, Cronin M, Dearden J, Gramatica P, Martin YC, Todeschini R, Consonni V, Kuz'min VE, Cramer R, Benigni R, Yang C, Rathman J, Terfloth L, Gasteiger J, Richard A, Tropsha A. QSAR modeling: where have you been? Where are you going to? J Med Chem. 2014 Jun 26;57(12):4977-5010. doi: 10.1021/jm4004285. Epub 2014 Jan 6. PMID: 24351051; PMCID: PMC4074254
- 9. Zhenyu Hou, Xiao Liu, Yukuo Cen, Yuxiao Dong, Hongxia Yang, Chunjie Wang, and Jie Tang. Graphmae: Self-supervised masked graph autoencoders. In Procs. of SIGKDD, 2022
- Aynacı Koyuncu, Elif & Yaşar, Ahmet & Arslan, Fatma & Sari, Nursen. (2019). Synthesis of novel Schiff base derivatives of tacrine and investigation of their acetylcholinesterase inhibition potency. Macedonian Journal of Chemistry and Chemical Engineering. 38. 75. 10.20450/mjcce.2019.1561