

Supplementary Information

Algebraic Graph-assisted Bidirectional Transformers for Molecular Prediction

Dong Chen^{1,2}, Kaifu Gao², Duc Duy Nguyen³, Xin Chen¹, Yi Jiang¹, Guo-Wei Wei ^{*2,4,5}
and Feng Pan ^{†1}

¹*School of Advanced Materials, Peking University, Shenzhen Graduate School, Shenzhen 518055, China*

²*Department of Mathematics, Michigan State University, MI, 48824, USA*

³*Department of Mathematics, University of Kentucky, KY 40506, USA*

⁴*Department of Electrical and Computer Engineering, Michigan State University, MI 48824, USA*

⁵*Department of Biochemistry and Molecular Biology, Michigan State University, MI 48824, USA*

This document contains supplementary information about methods and results which were not necessary to include in the central part of the paper but might be of interest to readers. This supplementary material includes the following sections:

Contents

S1 Datasets	2
S2 AGBT model parametrization	3
S2.1 Input processing	3
S2.2 Bidirectional transformer model parametrization	3
S2.3 Algebraic graph model parametrization	5
S2.4 Feature fusion	7
S2.5 Downstream machine learning algorithms	8
S3 Supplementary Figures	10
S4 Supplementary Tables	13

*Corresponding author: weig@msu.edu

†Corresponding author: panfeng@pkusz.edu.cn

S1 Datasets

In this work, we use ChEMBL [1] as the pre-trained dataset. ChEMBL is a chemical database of bioactive molecules with drug-like properties and it is a free database open to the public. [2] The ChEMBL26 is the current version of ChEMBL, updated in March 2020. There are over 1.9 million molecules in the ChEMBL26 dataset.

Table S1: The quantitative summary of all datasets.

	Datasets	Total size	Train set size	Test set size	Max value	Min value
Unlabeled data (pre-train)	ChEMBL	1936342	1926342	10000	-	-
Labeled data (fine-tune)	LD50	7413	5931	1482	7.201	0.291
	IGC50	1792	1434	358	6.36	0.334
	LC50	823	659	164	9.261	0.037
	LC50DM	353	283	70	10.064	0.117
	Log P	8605	8199	406	8.42	-4.64

For downstream tasks, four toxicity datasets were studied in our work, namely oral rate LD50, 40 h Tetrahymenapyriformis IGC50, 96 h fathead minnow LC50, and 48 h Daphnia Magna LC50DM, the basic information of toxicity datasets are shown in Table S1. Among them, LD50 measures the number of chemicals that can kill half of the rats when orally ingested. The LD50 represents the amount of chemicals that can kill half of the rats when orally ingested. It was originally from <https://chem.nlm.nih.gov/chemidplus/>. IGC50 records the 50% growth inhibitory concentration of Tetrahymena pyriformis organism after 40h. [3, 4] LC50 reports at the concentration of test chemicals in the water in milligrams per liter that cause 50% of fathead minnows to die after 96h. The last one is LC50DM, which represents the concentration of test chemicals in the water in milligrams per liter that cause 50% Daphnia Magna to die after 48h. LC50 and LC50DM were original from <http://cfpub.epa.gov/ecotox/>. The unit of toxicity reported in these four data sets is $-\log_{10}$ mol/L. The sizes of these four data sets vary from 353 to 7413, which poses a challenge for a predictive model to achieve consistent accuracy and robustness. For the partition coefficient prediction task, the training set contained 8199 molecules and the test set included 406 components. All components in the test set were approved as organic drugs by the Food and Drug Administration (FDA). The log P values for all training and test sets were compiled by Cheng et al. [5], and all log P values ranged from -4.64 to 8.42 (Table S1).

Table S2: A total of 51 symbols are used to split SMILES strings

Index	0	1	2	3	4	5	6	7	8	9
Symbol	c	C	()	O	1	2	=	N	@
Index	10	11	12	13	14	15	16	17	18	19
Symbol	[]	n	3	H	F	4	-	S	Cl
Index	20	21	22	23	24	25	26	27	28	29
Symbol	/	s	o	5	+	#	.	\	Br	6
Index	30	31	32	33	34	35	36	37	38	39
Symbol	P	I	7	Na	%	8	B	9	Si	0
Index	40	41	42	43	44	45	46	47	48	49
Symbol	Se	K	se	Li	As	Zn	Ca	Mg	Al	Te
Index	50									
Symbol	te									

Additionally, we statistic the length of SMILES for all molecules. As listed in Table S2, a total of 51

symbols are used to split these SMILES strings. The distribution of SMILES string lengths in the ChEMBL is shown in Figure S1a, and the majority of SMILES are within 254 in length. Therefore, in this work, we choose data with SMILES length no greater than 254 to pre-train. The exact number in the training set is 1,926, 342, and 10 thousand SMILES strings were randomly selected as a validating set. The basic information of ChEMBL used in pre-training is shown in Table S1. The distributions of SMILES string lengths for toxicity and logP datasets are shown in Figure S1b. Only one SMILES string on the LD50 dataset has a length of more than 254, with a length of 284. Therefore, in the downstream tasks, we truncate these sequences that exceeded the limit length and input only the first 254 symbols. All these datasets are also available at <https://weilab.math.msu.edu/Database/>.

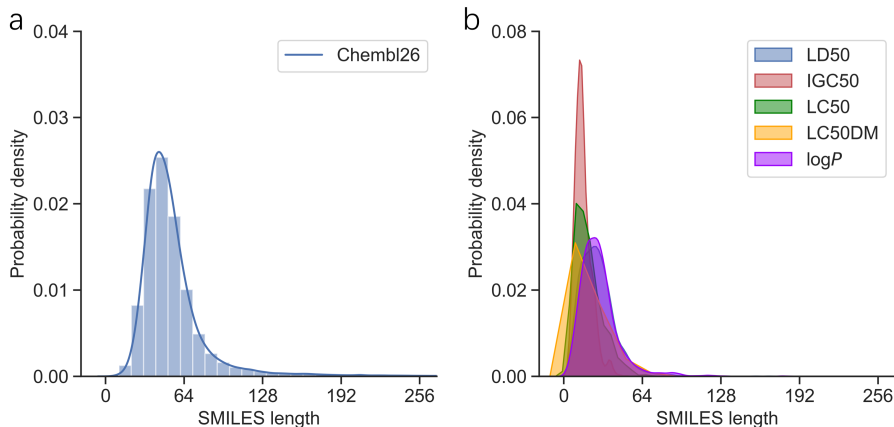


Figure S1: The distributions of SMILES string lengths. **a.** The ChEMBL database. **b.** Four toxicity datasets.

S2 AGBT model parametrization

S2.1 Input processing

In this work, all input SMILES strings for bidirectional transformer need to be processing. A total of 51 symbols, as listed in Table S2, are used to split these SMILES strings. We add a ' $\langle s \rangle$ ' symbol and a ' $\langle \backslash s \rangle$ ' at the beginning and end of each input SMILES, which represent the beginning and the end of each input, respectively. Besides, the ' $\langle unk \rangle$ ' is used to represent some undefined symbols. Since the length of SMILES varies from molecule to molecule, the ' $\langle pad \rangle$ ' is used as a padding symbol to fill in short inputs to reach the preset length. For the self-supervised learning (SSL) -based pre-training, the 15% symbol of the input SMILES needs to be operated. Among these 15% symbols, 80% of symbols were masked, 10% of the symbols were unchanged, and the remaining 10% were randomly replaced.

S2.2 Bidirectional transformer model parametrization

SSL-based pre-training Similar with the architecture of bidirectional encoder representations from transformers (BERT)[6], our pre-training model is a multi-layer bidirectional transformer encoder, as shown in Figure S2. Each transformer layer contains two sub-layers. The first is a multi-head self-attention layer, and the second is a fully connected feed-forward neural network. The residual connection is applied to each of the two sub-layers, followed by layer normalization.[7] Each transformer layer maps the output features from the former transformer layer or the embedded features from the input into different nonlinear space. The attention mechanism used in the transformer encoder is scaled dot-product attention and it is formulated as

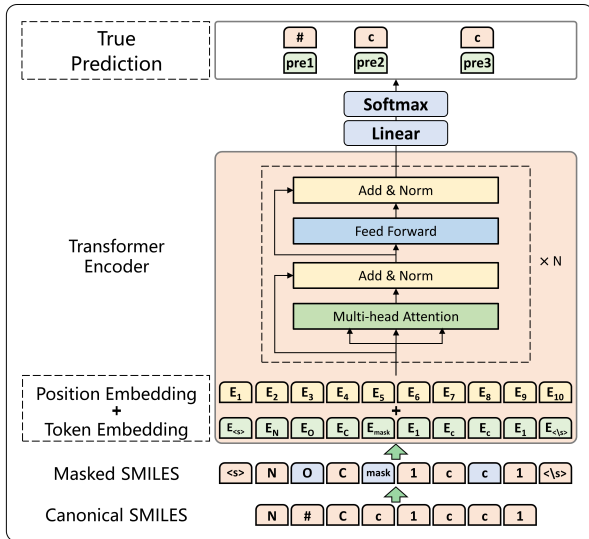


Figure S2: The whole structure of the bidirectional encoder from transformers used in pre-training.

follow,

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V. \quad (1)$$

The Q , K , and V , named query matrix, key matrix, and value matrix, are mapping from input data. The dot products of the query matrix and key matrix are divided by the scaling factor $\sqrt{d_k}$, where the d_k is the embedding dimension. In practice, a multi-head self-attention mechanism is applied in the transformer encoder, where different heads could pay attention to various aspects and improve performance. On the top of the N transformer encoder layers, there is a linear layer transforming the embedding dimension into the vocabulary size. Finally, the softmax function is used to select the maximum probability value of each masked location and report the corresponding prediction result.

The proposed model is based on Fairseq [8], which is a Sequence-to-Sequence Toolkit written in Python and PyTorch [9], and slightly modified so that it can be used for molecular analysis. In this work, the pre-training bidirectional transformer model contains 8 Transformer encoder layers, the embedding dimension is set to 512, the number of self-attention heads is 8, and the embedding size of fully connected feed-forward layers is 1024. The maximum sequence length is set to 256, including the start and terminate symbols. For better convergence, the Adam optimizer [10] is used in the pre-training and fine-tune, the Adam betas are (0.9, 0.999), and the weight decay is 0.1. Besides, a warming-up strategy is applied for the first 4000 updates and the total update steps are one million, the maximum learning rate is set to 0.0001 in this strategy. The cross-entropy was applied to measure the difference between the predicted symbols and the real symbols at the masked position. The model is trained on six Tesla V100-SXM2 GPUs and the maximum sequence number in each GPU is set to 64.

SSL-based and SL-based fine-tuning There are two strategies to be used in the fine-tuning stage: self-supervised learning (SSL) -based fine-tuning of task-specific data without using their labels and sequential supervised learning (SL) -based fine-tuning of task-specific data with their labels. For SSL-based fine-tune, the pre-trained model is fed with the input data of the downstream task-specific datasets, including both training sets and test sets. For each SMILES string, we randomly select 15% symbols to be a training-validation set in our loss function. Only 50% symbols of the set were masked and the remaining 50% symbols of the set were unchanged. A warming-up strategy is also applied for the first 500 updates. The total update steps are 2000. The maximum learning rate is set to 0.00001 in this stage. In the last hidden

layer, the embedded vector of length 512 correspondings to the first special symbol $\langle s \rangle$ is used for molecular property prediction. Figure S3 shows the workflow of generating molecular fingerprints from the SSL-based fine-tuned model.

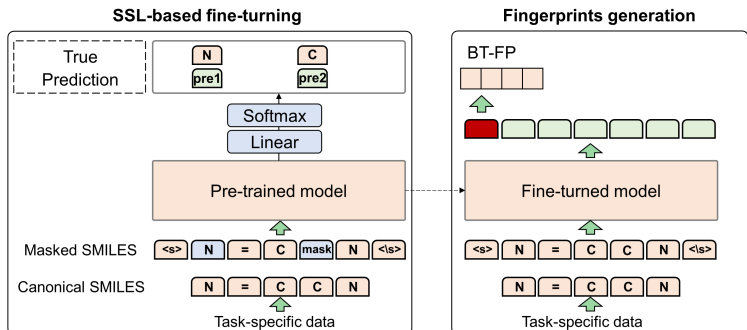


Figure S3: Workflow for generating molecular fingerprints from the pre-trained and SSL-based fine-tuned model. Three two mask operations, 'mask' and 'no changing', are retained in the self-supervised fine-tuning stage. The labels of task-specific data are disregarded in the SSL-based fine-tuning stage. Here, $\langle s \rangle$ is a special leading symbol added in front of every input SMILES, and $\langle \backslash s \rangle$ is a terminating symbol. At the stage of fingerprints generation, $\langle s \rangle$'s embedding vector from the bidirectional encoder is utilized to represent the molecular fingerprint (BT-FP).

For sequential SL-based fine-tuning, the labels of task-specific data are utilized. The pre-trained model will be fed with data from the training set of the task-specific dataset, and no additional 'mask' operations are required for the input SMILES. The Adam optimizer is set as the same as that of pre-training. The maximum learning rate is set to 10^{-5} . The warming-up strategy is used for the first 500 updates and the total update steps are 5000 for each dataset. The mean square error metric is used in this fine-tuning stage, as shown in Figure S4. All models were trained on six Tesla V100-SXM2 GPUs and the maximum sequence number in each GPU is set to 64.

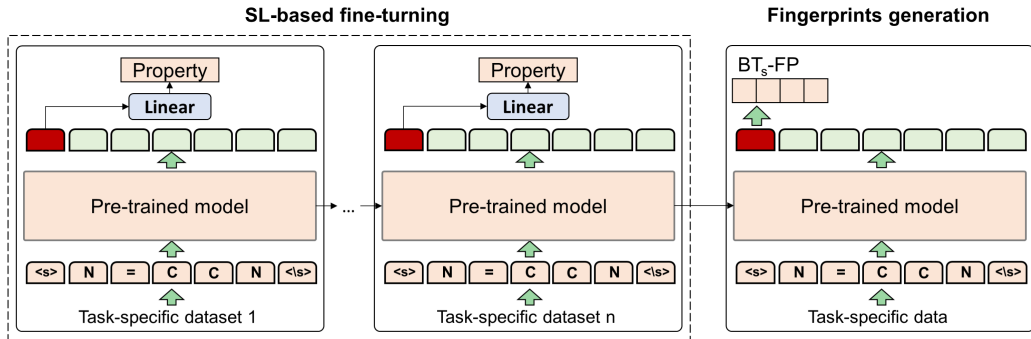


Figure S4: Workflow for generating molecular fingerprints from the pre-trained and sequential supervised learning (SL) fine-tuned model. Labeled task-specific data are employed in the sequential SL-based fine-tuning stage. At the stage of fingerprints generation, $\langle s \rangle$'s embedding vector before last linear layer is utilized to represent the molecular fingerprint (BT_s-FP).

S2.3 Algebraic graph model parametrization

In order to select a general AG-FP for all four toxicity data sets, we need to combine the kernel function and graph matrix type properly. For the sake of convenience, we use the notation $AG_{\Omega, \beta, \tau}^{\mathcal{M}}$, to indicate the AG-FPs generated by using interactive matrix type \mathcal{M} with kernel function ω and corresponding kernel parameters β and τ . Here, $\mathcal{M} = \{Adj, Lap\}$ represents a set of adjacency matrix and Laplacian matrix.

$\omega = \{E, L\}$ refers to a set of generalized exponential and generalized Lorentz kernels. In addition, the kernel parameter $\beta = \kappa$ if $\omega = E$, and $\beta = v$ if $\omega = L$. And τ is used such that $\eta_{k_1 k_2} = \tau(\bar{r}_{k_1} + \bar{r}_{k_2})$, where \bar{r}_{k_1} and \bar{r}_{k_2} are the van der Waals radii of element type k_1 and k_2 , respectively. Kernel parameters β and τ as selected based on the cross validation with a random split of the training data. It has been shown that multiscale information can boost the performance of predictor. [11, 12] In this work, we consider at most two kernels. As a straightforward notation extension, two kernels can be parametrized by $\text{AG}_{\omega_1, \beta_1, \tau_1; \Omega_2, \beta_2, \tau_2}^{\mathcal{M}_1, \mathcal{M}_2}$. To attain the best performance using AG-FP, the kernel parameters need to be optimized. We vary β , both τ and κ , from 0.5 to 6 with an increment of 0.5, while τ values are chosen from 0.5 to 6 with an increment of 0.5. The high values of the power order such as $\beta \in \{10, 15, 20\}$ are also considered to approximate the idea low-pass filter.[13] We use the method of 5-fold cross-validation (CV) to select the kernel hyperparameters \mathcal{M} , Ω , β and τ . Figure S5a shows the CV results of the single-kernel model ($\text{AG}_{\omega_1, \beta_1, \tau_1}^{\mathcal{M}_1}$), and R^2 is used as the evaluation metrics. Then based on the optimal kernel parameters in the single-kernel model, the two-kernel model, $\text{AG}_{\omega_1, \beta_1, \tau_1; \Omega_2, \beta_2, \tau_2}^{\mathcal{M}_1, \mathcal{M}_2}$, can be optimized by using 5-fold CV on training sets. Figure S5b in the supplement material reports the best models with associated R^2 in this experiment. All cross-validations were performed for toxicity training sets, and the scores were based on the mean value of R^2 in the four training sets.

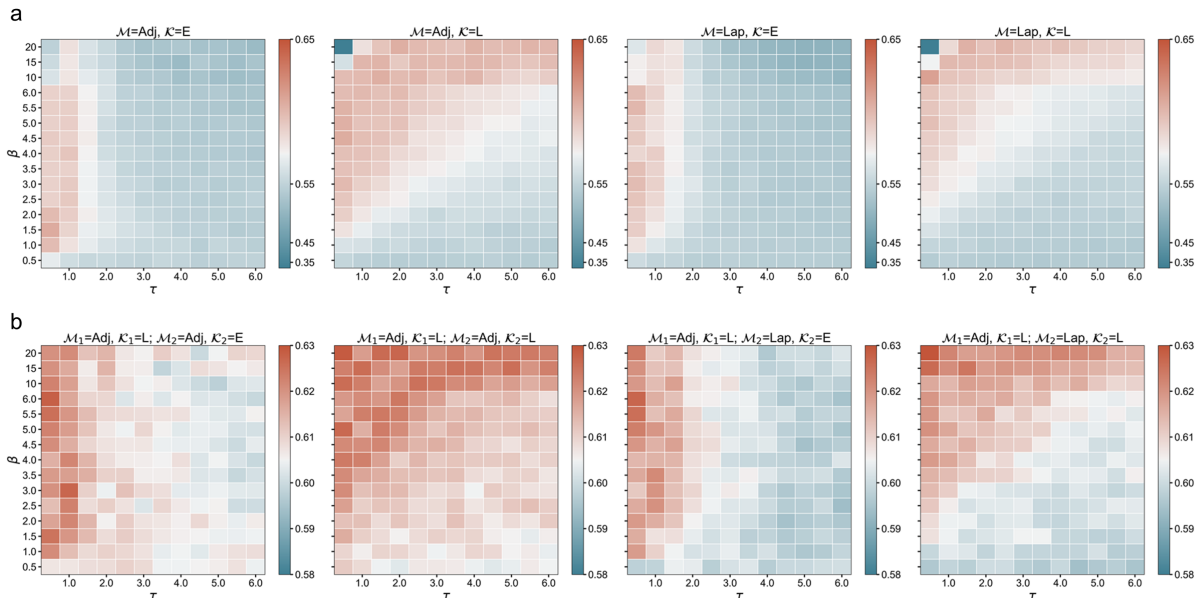


Figure S5: Squared Pearson correlation coefficients (R^2) from 5-fold cross-validation of $\text{AG}_{\Omega, \beta, \tau}^{\mathcal{M}}$, and $\text{AG}_{\omega_1, \beta_1, \tau_1; \Omega_2, \beta_2, \tau_2}^{\mathcal{M}_1, \mathcal{M}_2}$ on the training data of four toxicity datasets are plotted against different values of τ and β . **a.** The best hyperparameters and R^2 for these one-scale models are found to be ($\text{AG}_{E, 1.5, 0.5}^{\text{Adj}}$, average $R^2 = 0.616$), ($\text{AG}_{L, 4.5, 0.5}^{\text{Adj}}$, average $R^2 = 0.616$), ($\text{AG}_{E, 5.5, 0.5}^{\text{Lap}}$, average $R^2 = 0.610$) and ($\text{AG}_{L, 10, 0.5}^{\text{Lap}}$, average $R^2 = 0.620$) from left to right separately. **b.** Based on the best one-scale model, the best hyperparameters and R^2 for these multiscale models are found to be ($\text{AG}_{L, 10, 0.5; E, 6, 0.5}^{\text{Lap, Adj}}$, average $R^2 = 0.628$), ($\text{AG}_{L, 10, 0.5; L, 20, 0.5}^{\text{Lap, Adj}}$, average $R^2 = 0.629$), ($\text{AG}_{L, 10, 0.5; E, 6, 0.5}^{\text{Lap, Lap}}$, average $R^2 = 0.627$) and ($\text{AG}_{L, 10, 0.5; E, 20, 0.5}^{\text{Lap, Lap}}$, average $R^2 = 0.631$) from left to right separately.

For the toxicity and logP datasets, there are 10 commonly occurring element types, i.e., {H, C, N, O, F, P, S, Cl, Br, I}, which means 100 element interactive pairs will form based on the combinations of these 10 element types in molecules. For adjacency matrices, only positive eigenvalues are considered. Note that Laplacian matrices are positive semidefinite. As discussed in the Methods section, we can compute nine descriptive statistical values, namely the maximum, minimum, average, summation, median, standard deviation, and variance of all eigenvalues. Another two values are the number of considered eigenvalues and the sum of the second power of eigenvalues. This gives rise to a total of 900 features for one kernel,

which means that we can get an 1800 dimension AG-FP for each molecule if we use two-kernel information. The optimal two-kernel algebraic graph models are $AG_{L,10,0.5;L,20,0.5}^{Lap,Lap}$, and the average R^2 of all four toxicity datasets is 0.631.

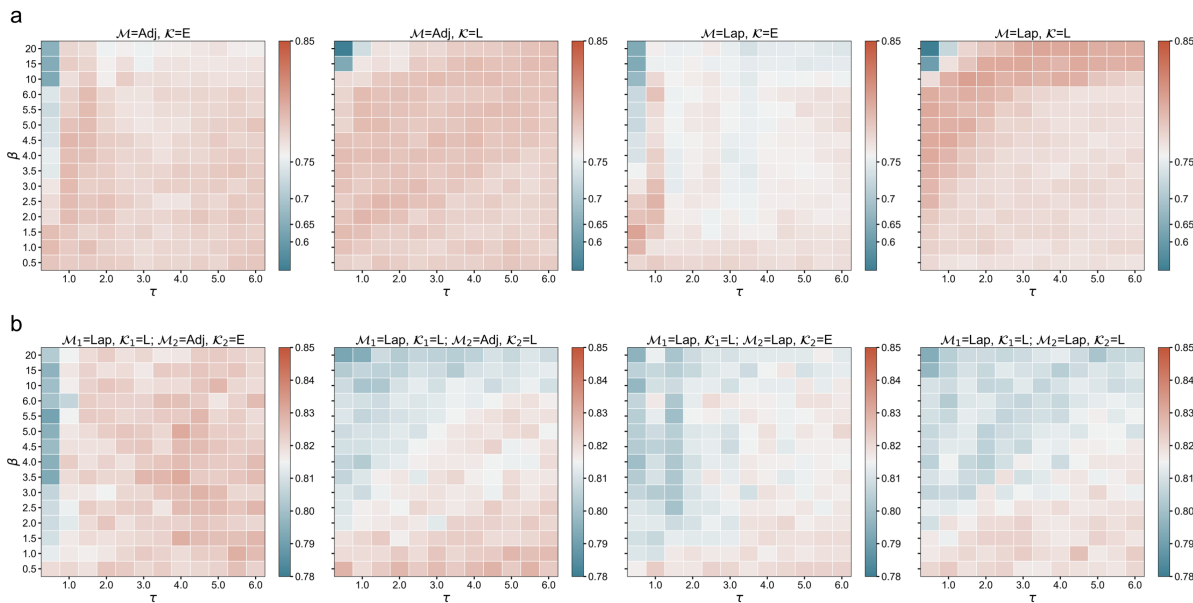


Figure S6: Squared Pearson correlation coefficients (R^2) from 5-fold cross-validation of $AG_{\Omega,\beta,\tau}^{\mathcal{M}}$, and $AG_{\omega_1,\beta_1,\tau_1;\omega_2,\beta_2,\tau_2}^{\mathcal{M}_1,\mathcal{M}_2}$ on the training data of partition coefficient data sets are plotted against different values of τ and β . **a.** The best hyperparameters and R^2 for these one-scale models are found to be ($AG_{E,4.5,1.0}^{Adj}$, $R^2 = 0.798$), ($AG_{L,2.0,1.0}^{Adj}$, $R^2 = 0.799$), ($AG_{E,1.5,1.0}^{Lap}$, $R^2 = 0.81$) and ($AG_{L,10,1.5}^{Lap}$, $R^2 = 0.811$) from left to right separately. **b.** Based on the best one-scale model, the best hyperparameters and R^2 for these multiscale models are found to be ($AG_{L,10,1.5;E,5,4}^{Lap,Adj}$, $R^2 = 0.831$), ($AG_{L,10,1.5;L,0.5,5.5}^{Lap,Adj}$, $R^2 = 0.829$), ($AG_{L,10,1.5;E,0.5,1}^{Lap,Lap}$, $R^2 = 0.823$) and ($AG_{L,10,1.5;E,1,4.5}^{Lap,Lap}$, $R^2 = 0.826$) from left to right separately.

For the partition coefficient dataset, we also use two-kernel information as the final AG-FPs. There are the same 10 element types as for toxicity datasets and there are also 100 element interactive pairs will form in each molecule. As shown in Figure S6a, we firstly selected the best hyperparameters for these one-scale models, which are $AG_{E,4.5,1.0}^{Adj}$ ($R^2 = 0.798$), $AG_{L,2.0,1.0}^{Adj}$ ($R^2 = 0.799$), $AG_{E,1.5,1.0}^{Lap}$ ($R^2 = 0.81$), and $AG_{L,10,1.5}^{Lap}$ ($R^2 = 0.811$). Based on the optimal one-scale model ($AG_{L,10,1.5}^{Lap}$), the best multiscale model is found to be $AG_{L,10,1.5;E,5,4}^{Lap,Adj}$, as shown in Figure S6b, the value of R^2 is 0.831. The gradient boosting decision tree (GBDT) is used to select optimal algebraic graph model hyperparameters. The parameters of GBDT vary with the size of the training set, which are listed in Table S3.

S2.4 Feature fusion

Based on a large amount of unlabeled data, BT-FP can capture the overall information of molecules after pre-training and fine-tuning. AG-FP, on the other hand, as insight based on physical and chemical knowledge, can obtain more detailed information of molecular structure, including dihedral angle and relative distance of atoms, with the help of algebraic graph theory. The proposed AGBT-FP in this work is a fusion of BT-FP and AG-FP. The random forest (RF) is used to fuse BT-FP and AG-FP. First, we combine BT-FP and AG-FP. Then the RF algorithm was used to select top 512 features. The parameters of RF vary with the size of the training set. All parameters are listed in Table S3. The final AGBT-FP's dimension is set to 512, which is the same as BT-FP's.

S2.5 Downstream machine learning algorithms

To compare the AGBT and other fingerprints’ performance on specific tasks, three machine learning algorithms are used: gradient boosting decision tree (GBDT), single-task deep neural network (ST-DNN), and multitask deep neural network (MT-DNN).

Gradient boosting decision tree (GBDT). GBDT is a robust machine learning regressor. In this approach, individual decision trees are successively combined in a stage-wise fashion to achieve the capability of learning complex features. It uses both gradient and boosting strategies to reduce model errors. Compared to the deep neural network (DNN) approaches, this ensemble method is robust, relatively insensitive to hyperparameters, and easy to implement. Moreover, they are much faster to train than DNN. In fact, for small data sets, GBDT can perform even better than DNN or other deep learning algorithms.[14, 15] Therefore, GBDT has been applied to a variety of QSAR prediction problems, such as toxicity, solvation, and binding affinity predictions.[16, 17]

The GBDT is used to predict the toxicity and $\log P$ in this work and implemented by the scikit-learn package.[18] In this work, there are five data sets with their training data size varying from 283 to 8199. To better compare feature performance, we set only two sets of parameters according to the size of the training set for GBDT. The detailed values of these hyperparameters are given in Table S3.

Table S3: RF and GBDT parameters for different toxicity training-set sizes

Training-set Size	RF Parameters	GBDT Parameters
> 1000	n_estimators = 10000 criterion = ‘mse’ max_depth = 8 min_samples_split = 4 min_samples_leaf = 1 min_weight_fraction_leaf = 0.0	n_estimators = 10000 max_depth = 8 min_samples_split = 4 learning_rate = 0.01 subsample = 0.3 max_features=‘sqrt’
< 1000	n_estimators = 10000 criterion = ‘mse’ max_depth = 7 min_samples_split = 3 min_samples_leaf = 1 min_weight_fraction_leaf = 0.0	n_estimators = 10000 max_depth = 7 min_samples_split = 3 learning_rate = 0.01 subsample = 0.2 max_features=‘sqrt’

Single-task deep neural network (ST-DNN). A DNN mimics the learning process of a biological brain by constructing a wide and deep architecture of numerous connected neuron units. A typical deep neural network often includes multiple hidden layers. In each layer, there are hundreds or even thousands of neurons. During learning stages, weights on each layer are updated by backpropagation. With a complex and deep network, DNN is capable of constructing hierarchical features and model complex nonlinear relationships. ST-DNN is a regular deep learning algorithm. It only takes care of one single prediction task. Therefore, it only learns from one specific training dataset. A typical four-layer ST-DNN is showed in Figure S7a, where N_i ($i = 1, \dots, 4$), represents the number of neurons in the i th hidden layer.

Multitask deep neural network (MT-DNN). The multitask (MT) learning technique has achieved much success in qualitative Merck and Tox21 prediction challenges.[19, 20, 21] In the MT framework, multiple tasks share the same hidden layers. However, the output layer is attached to different tasks. This framework enables the neural network to learn all the data simultaneously for different tasks. Thus, the commonalities and differences among various data sets can be exploited. It has been shown that MT learning typically can improve the prediction accuracy of relatively small data sets if it combines with relatively larger data sets in its training. Figure S7b is an illustration of a typical four-layer MT-DNN for training four different tasks simultaneously. Suppose there are totally T tasks and the training data for the t th task are $(X_i^t, y_i^t)_{i=1}^{N_t}$,

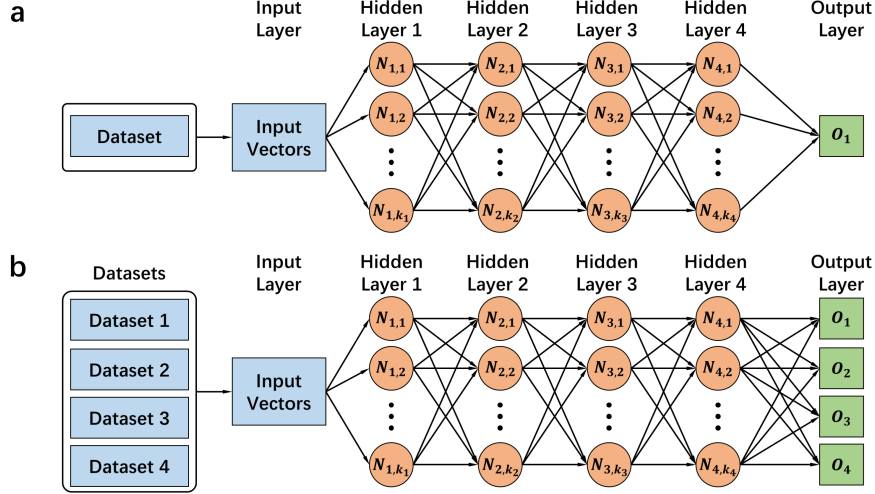


Figure S7: ST-DNN and MT-DNN framework. a) An illustration of a typical ST-DNN. Only one dataset is trained in this network. Four hidden layers are included, k_i ($i = 1, 2, 3, 4$) represents the number of neurons in the i th hidden layer and $N_{i,j}$ is the j th neuron in the i th hidden layer. Here, O_1 is the single output for the model. b) An illustration of a typical MT-DNN training four tasks (datasets) simultaneously. Four hidden layers are included in this network, k_i ($i = 1, 2, 3, 4$) represents the number of neurons in the i th hidden layer and $N_{i,j}$ is the j th neuron in the i th hidden layer. Here O_1 to O_4 represent four predictor outputs for four tasks.

where $t = 1, \dots, T$, $i = 1, \dots, N_t$, where N_t is the number of samples in the t th task, and X_i^t is the feature vector for the i th sample in the t th task, y_i^t is the label value of the i th sample in the t th task, respectively. The purpose of MT learning is to simultaneously minimize the loss function:

$$\operatorname{argmin} \sum_{t=1}^T \sum_{i=1}^{N_t} L(y_i^t, f^t(X_i^t, \theta^t)), \quad (2)$$

where f^t is the prediction for the i th sample in the t th task by our MT-DNN, which is a function of the feature vector X_i^t , L is the loss function, and θ^t is the collection of machine learning hyperparameters. A popular cost function for regression is the mean squared error, which is formulated as:

$$L(y_i^t, f^t(X_i^t, \theta^t)) = \frac{1}{N_t} \sum_{i=1}^{N_t} (y_i^t - f^t(X_i^t, \theta^t))^2. \quad (3)$$

In this work, MT-DNN is only applied to predict the toxicity. The ultimate goal of MT-DNN learning is to potentially improve the overall performance of multiple toxicity prediction models, especially for the smallest dataset that performs relatively poorly in the ST-DNN. More concretely, it is reasonable to assume that different toxicity indices share a common statistic pattern so that these different tasks can be trained simultaneously when their feature vectors are constructed in the same manner. For our toxicity prediction, four different tasks (LD50, IGC50, LC50, LC50DM data sets) are trained together. This leads to four output neurons in the output layer, with each neuron being specific to one of four tasks.

The performance of deep neural network models depends on their architecture, input data dimension, and hyperparameters. For BT-FP and AGBT-FP, the feature sizes are both 512, which means that the network with the same architecture can be used to train these two sets of features. The input layer contains 512 neurons, followed by four hidden layers with 1024, 512, 512, and 512 neurons, respectively. For the present regression problem, only one neuron in the final output layer. For AG-FP, it contains 1800 features, and thus a more complex network structure is required. In this case, we set 1800 neurons in the input layer, followed by 5 hidden layers with 2048, 1024, 512, 512, and 512 neurons, respectively. The output layer has one neuron. Other network parameters are all the same for these three kinds of molecular features. The

211 stochastic gradient descent (SGD) with a momentum of 0.5 is used as an optimizer. We use 2000 epochs
 212 to train all the networks. The mini-batch size is set to 8. The learning rate is set to 0.01 in the first 1000
 213 epochs and 0.001 for the rest epochs. These hyperparameters are applied to both ST-DNN and MT-DNN.
 214 All the DNN models are built and trained in Pytorch.[9]

215 S3 Supplementary Figures

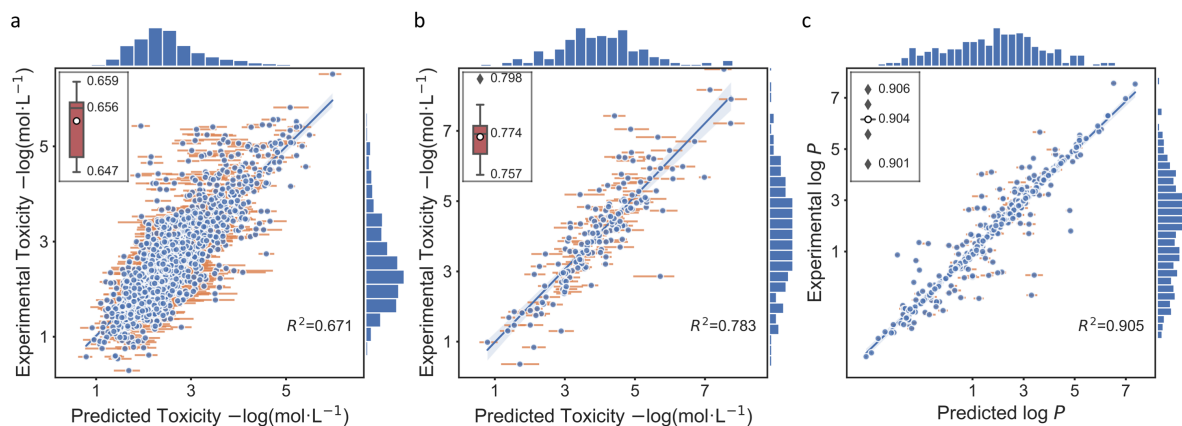


Figure S8: Data and results of AGBT. **a**, Predicted results of AGBT-FPs with MT-DNN model for LD50 dataset ($R^2=0.671$, RMSE=0.554 log(mol/L) for 20 repeated experiments). **b**, Predicted results of BT-FPs with MT-DNN model for LC50 dataset ($R^2=0.783$, RMSE=0.692 log(mol/L) for 20 repeated experiments). **c**, Predicted results of AGBT_s-FPs with MT-DNN model for LC50 dataset ($R^2=0.905$, RMSE=0.615 log(mol/L) for 20 repeated experiments).

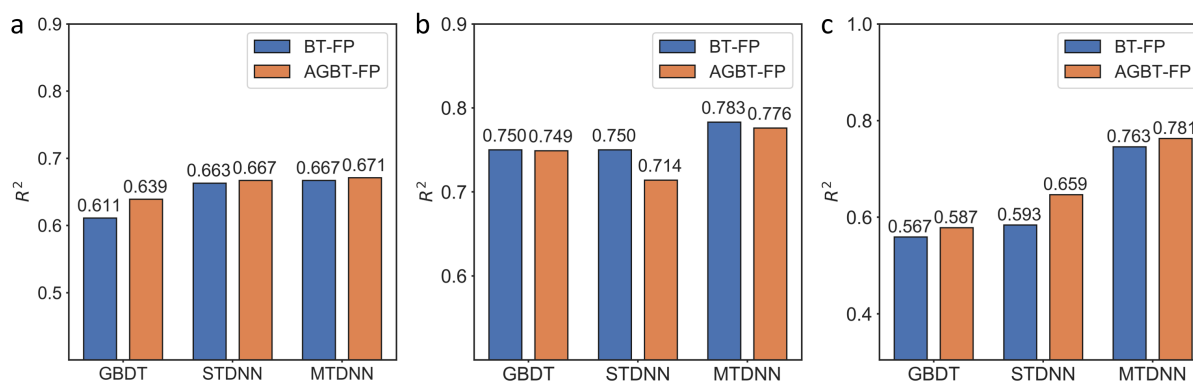


Figure S9: R^2 values of BT-FP and AGBT-FP predictions on three machine learning algorithms, GBDT, STDNN, and MTDNN. **a** LD50 dataset, **b** LC50 dataset, **c** LC50DM dataset.

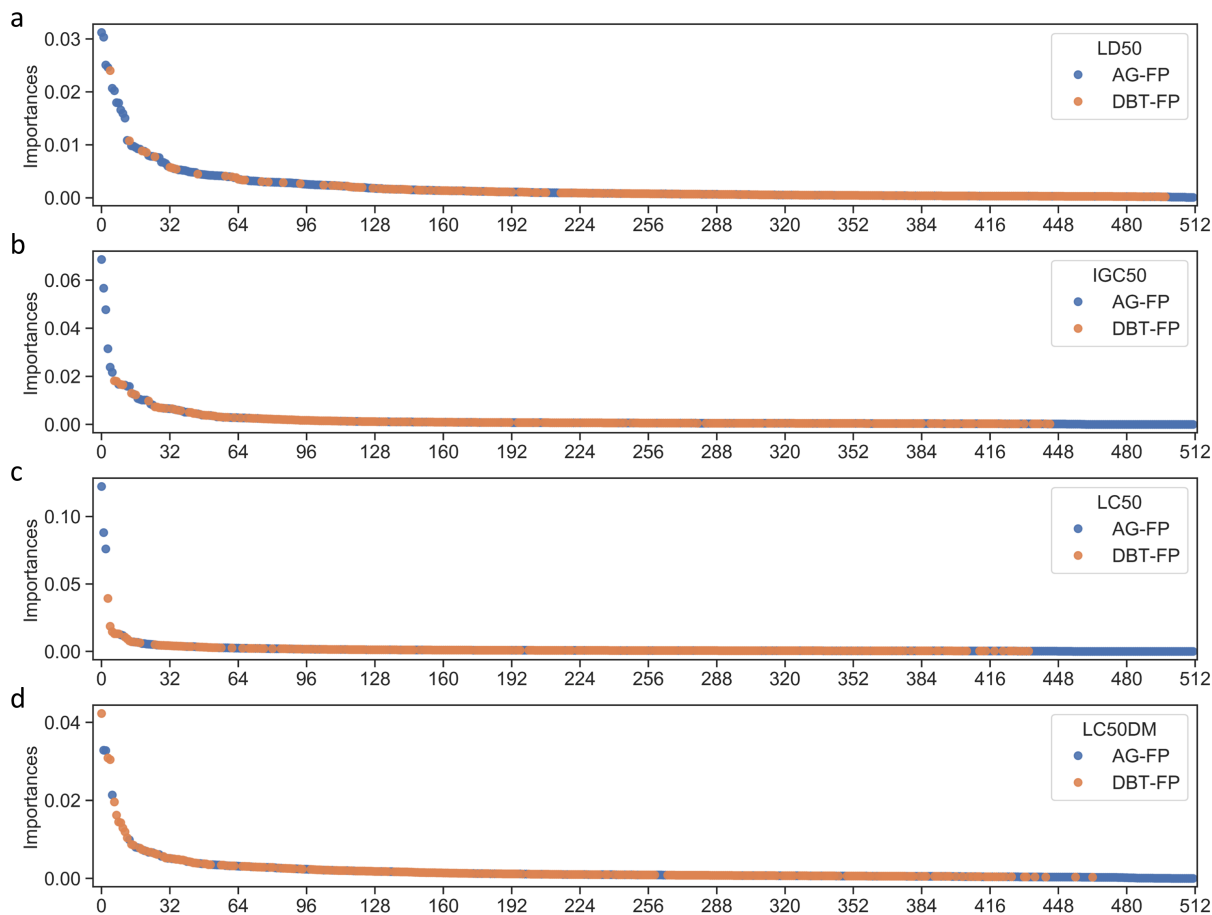


Figure S10: The AGBT-FPs of the four toxicity datasets were ranked by their feature importance. **a** Sorted feature importance for the LD50 dataset. The top three features are from AG-FP. **b** Sorted feature importance for the IGC50 dataset. The top three features are from AG-FP. **c** Sorted feature importance for the LC50 dataset. The top three features are from AG-FP. **d** Sorted feature importance for the LC50 dataset. The most important feature is from BT-FP. The 2nd and 3rd most important features are from AG-FP.

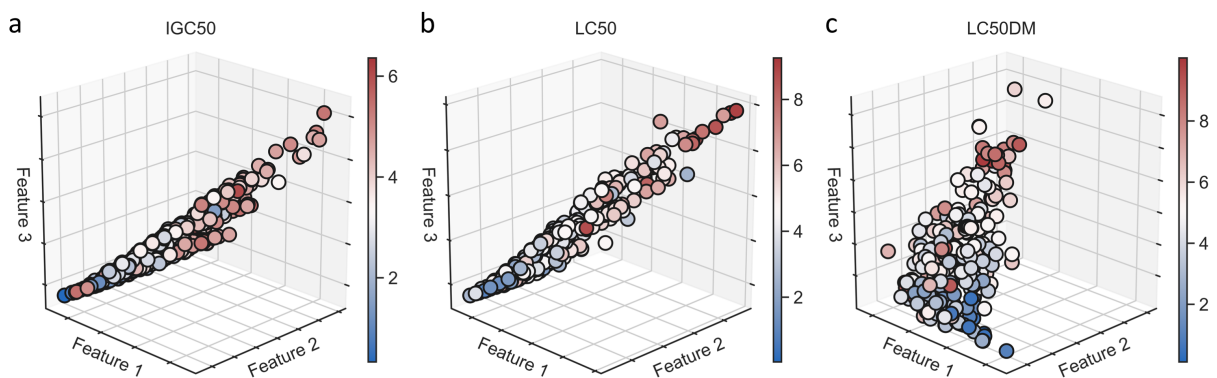


Figure S11: Distribution of molecules in the three most important features of AGBT-FP. **a** The distribution of the IGC50 dataset. **b** The distribution of the LC50 dataset. **c** The distribution of the LC50DM dataset.

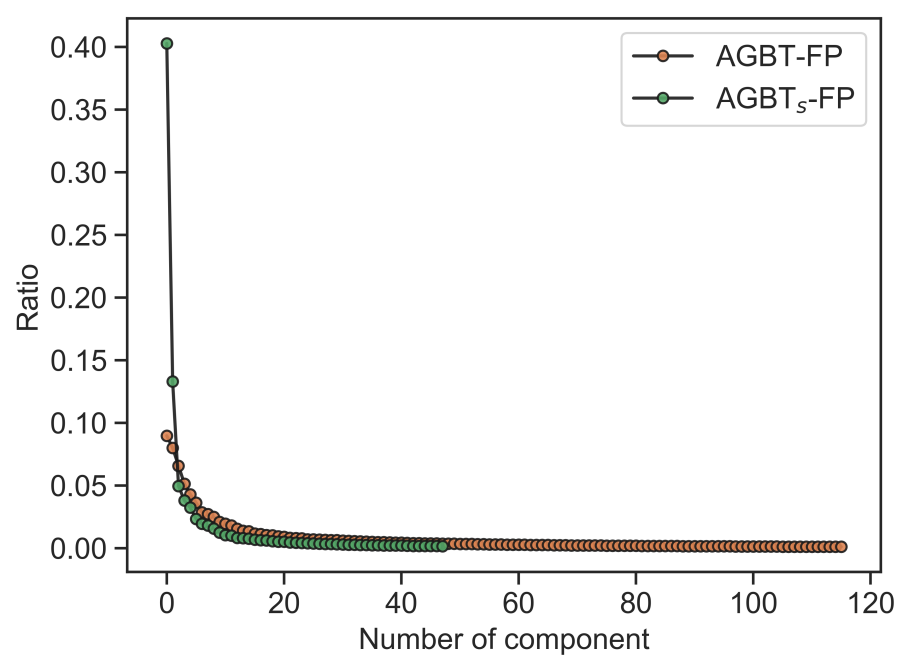


Figure S12: The ratio refers to the rate of variability (variance) of the data explained by each principal component through principal component analysis (PCA). For AGBT-FP (orange), the first 112 components are needed to represent 90% variance, whereas for AGBT_s-FP (green), only the first 48 components are needed to represent 90% of the variance.

S4 Supplementary Tables

Tanimoto coefficient, $S_{A,B}$, is used in this work to calculate the degree of similarity between two molecules. A higher average $S_{A,B}$ of the two datasets implies a higher similarity. Tanimoto coefficient, $S_{A,B}$, is defined as follow:

$$S_{A,B} = \frac{\sum_{i=1}^N x_{iA}x_{iB}}{\sum_{i=1}^N x_{iA}^2 + \sum_{i=1}^N x_{iB}^2 - \sum_{i=1}^N x_{iA}x_{iB}}. \quad (4)$$

In this study, the similarity between the largest dataset LD50, which contains 7413 molecules, with other three datasets are list in [Table S4](#)

Table S4: Similarity between the Largest Dataset LD50 with the other three datasets^a

Fingerprints	IGC50(1792)	LC50(823)	LC50DM(353)
Estate2	0.964	0.973	0.982
FP2	0.886	0.928	0.941

^a The number in the bracket is the total size of the dataset.

References

- [1] Anna Gaulton, Anne Hersey, Michal Nowotka, A Patricia Bento, Jon Chambers, David Mendez, Prudence Mutowo, Francis Atkinson, Louisa J Bellis, Elena Cibrián-Uhalte, et al. The chembl database in 2017. *Nucleic acids research*, 45(D1):D945–D954, 2017.
- [2] Anna Gaulton, Louisa J Bellis, A Patricia Bento, Jon Chambers, Mark Davies, Anne Hersey, Yvonne Light, Shaun McGlinchey, David Michalovich, Bissan Al-Lazikani, et al. ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic acids research*, 40(D1):D1100–D1107, 2012.
- [3] Kevin S Akers, Glendon D Sinks, and T Wayne Schultz. Structure–toxicity relationships for selected halogenated aliphatic chemicals. *Environmental toxicology and pharmacology*, 7(1):33–39, 1999.
- [4] Hao Zhu, Alexander Tropsha, Denis Fourches, Alexandre Varnek, Ester Papa, Paola Gramatica, Tomas Oberg, Phuong Dao, Artem Cherkasov, and Igor V Tetko. Combinatorial qsar modeling of chemical toxicants tested against tetrahymena pyriformis. *Journal of chemical information and modeling*, 48(4):766–784, 2008.
- [5] Tiejun Cheng, Yuan Zhao, Xun Li, Fu Lin, Yong Xu, Xinglong Zhang, Yan Li, Renxiao Wang, and Luhua Lai. Computation of octanol- water partition coefficients by guiding an additive model with knowledge. *Journal of chemical information and modeling*, 47(6):2140–2148, 2007.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [7] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [8] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. *arXiv preprint arXiv:1904.01038*, 2019.

- [9] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [10] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [11] Kristopher Opron, Kelin Xia, and Guo-Wei Wei. Communication: Capturing protein multiscale thermal fluctuations, 2015.
- [12] Duc D Nguyen, Tian Xiao, Menglun Wang, and Guo-Wei Wei. Rigidity strengthening: A mechanism for protein–ligand binding. *Journal of chemical information and modeling*, 57(7):1715–1721, 2017.
- [13] Kelin Xia, Kristopher Opron, and Guo-Wei Wei. Multiscale gaussian network model (mgnm) and multiscale anisotropic network model (manm). *The Journal of chemical physics*, 143(20):11B616-1, 2015.
- [14] Kaifu Gao, Duc Duy Nguyen, Vishnu Sresht, Alan M Mathiowetz, Meihua Tu, and Guo-Wei Wei. Are 2d fingerprints still valuable for drug discovery? *Physical Chemistry Chemical Physics*, 22(16):8373–8390, 2020.
- [15] Jian Jiang, Rui Wang, Menglun Wang, Kaifu Gao, Duc Duy Nguyen, and Guo-Wei Wei. Boosting tree-assisted multitask deep learning for small scientific datasets. *Journal of Chemical Information and Modeling*, 60(3):1235–1244, 2020.
- [16] T Martin et al. User’s guide for test (version 4.2)(toxicity estimation software tool): A program to estimate toxicity from molecular structure. *Washington (USA): US-EPA*, 2016.
- [17] Bao Wang, Chengzhang Wang, Kedi Wu, and Guo-Wei Wei. Breaking the polar-nonpolar division in solvation free energy prediction. *Journal of computational chemistry*, 39(4):217–233, 2018.
- [18] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- [19] Stephen J Capuzzi, Regina Politi, Olexandr Isayev, Sherif Farag, and Alexander Tropsha. Qsar modeling of tox21 challenge stress response and nuclear receptor signaling toxicity assays. *Frontiers in Environmental Science*, 4:3, 2016.
- [20] Bharath Ramsundar, Bowen Liu, Zhenqin Wu, Andreas Verras, Matthew Tudor, Robert P Sheridan, and Vijay Pande. Is multitask deep learning practical for pharma? *Journal of chemical information and modeling*, 57(8):2068–2076, 2017.
- [21] Jan Wenzel, Hans Matter, and Friedemann Schmidt. Predictive multitask deep neural network models for adme-tox properties: learning from large data sets. *Journal of chemical information and modeling*, 59(3):1253–1268, 2019.