

## A Dataset Details

Our study utilizes 15 molecule datasets of varying sizes obtained from MoleculeNet [14, 7] and ChemBL [12], which can be categorized into three groups based on their applications or source domains: medication (HIV, SIDER, Lipo, ClinTox), quantum mechanics (qm8, qm9), and chemical analysis (ESOL, FreeSolv, ChemBL, MUV, BACE, BBBP, ToxCast, Tox21, PCBA). As for qm8 and qm9, we randomly sample 3,000 graphs to construct the datasets. We use the original split setting, where qm8 and qm9 are randomly split, and scaffold splitting is used for the others. The small molecule datasets with less than 10,000 instances in the training sets are selected as target molecule datasets, i.e., Tox21, BACE, BBBP, ClinTox, SIDER, ToxCast, ESOL, FreeSolv, Lipo, qm8 and qm9. All the involved datasets can be accessed and downloaded from OGB<sup>1</sup> or MoleculeNet repository<sup>2</sup>. The overall statistics are summarized as follows:

	HIV	PCBA	Tox21	BACE	BBBP	ClinTox	MUV
#graphs	41,127	437,929	7,831	1,513	2,039	1,477	93,087
#tasks	1	128	12	1	1	2	17
split	scaffold	scaffold	scaffold	scaffold	scaffold	scaffold	scaffold
metric	ROAUC	AP	ROAUC	ROAUC	ROAUC	ROAUC	AP

Table 1: Dataset statistics(1).

	SIDER	ToxCast	ChemBL	ESOL	FreeSolv	Lipo	qm8	qm9
#graphs	1,427	8,576	456,309	1,128	642	4,200	3,000	3,000
#tasks	27	617	1,310	1	1	1	12	12
split	scaffold	scaffold	scaffold	scaffold	scaffold	scaffold	random	random
metric	ROAUC	ROAUC	ROAUC	RMSE	RMSE	RMSE	MAE	MAE

Table 2: Dataset statistics(2).

## B Implementation Details

The core code is provided in the supplementary material.

**Backbone model settings.** As for GIN [15], we fix the batch size as 128 and train the model for 50 epochs. We use Adam [9] with a learning rate of 0.001 for optimization. The hidden size and number of layers are set as 300 and 5 respectively. We set the dropout rate as 0.5 and apply batchnorm [8] in each layer. All the results are reported after 5 different random seeds.

As for Graphormer [16], we fix the batch size as 128 and train the model for 30 epochs. AdamW [11] with a learning rate of 0.0001 is used as the optimizer. The hidden size, number of layers, and number of attention heads are set as 512, 5, and 8 respectively. We set the dropout rate and attention dropout rate as 0.1 and 0.3. Layernorm [2] is applied across layers. The maximum number of nodes and the distance between nodes in the sampled graph is set as 128 and 5 respectively. The size of position embedding, in-degree embedding, and out-degree embedding are fixed as 512. All the results are reported after 5 different random seeds.

**Grouping method settings.** As for MolGroup, we apply GIN as the encoder. During the training, we randomly select datasets with equal probability and sample data from them to ensure that each dataset contributes an equal number of samples to the constructed mini-batch. We use MACCS [4] fingerprint to calculate the structure affinity score. We fix the number of the filtering iterations as 3 and the balance coefficient  $\lambda$  is set as 0.9. During each round, the auxiliary datasets with affinity scores below 0.6 will be filtered out. If none of the datasets fall below this threshold, we filter out the dataset with the lowest affinity score. The orthogonal initialization [13] is applied to initialize the learnable task embedding  $e^{\text{task}}$  to make sure that the dot product between the target dataset

<sup>1</sup><https://ogb.stanford.edu/>

<sup>2</sup><https://moleculenet.org/>

---

**Algorithm 1:** Iterative filtering process with MolGroup

---

**Input:** Target dataset  $\mathcal{D}_T$  with  $N$  training instances; Candidate auxiliary datasets  $\{\mathcal{D}_A\}_M$ ; GNN model with specific parameters  $\theta_T, \{\theta_m\}_M$  for each dataset and routing mechanism  $g(\cdot)$ ; Number of iterations  $R$ ; Number of epochs  $E$ ; Learning rate  $lr$ ; Batch size  $B$ .

```
// Filtering round
for  $r \leftarrow 1, \dots, R$  do
  Random initialize  $\theta_T, \theta_1, \dots, \theta_M$  and  $g(\cdot)$ .
   $\phi_1, \dots, \phi_M \leftarrow 0$ . // Affinity scores
   $I \leftarrow NM/B$ . // Number of iteration in each epoch
  // Training epoch
  for  $e \leftarrow 1, \dots, E$  do
    // Training step
    for  $iter \leftarrow 1, \dots, I$  do
      Sample mini-batch  $\{\mathcal{B}_T, \mathcal{B}_1, \dots, \mathcal{B}_M\}$  from current datasets.
      Obtain losses  $l_T, \{l_m\}_M$  and affinity scores  $\{\alpha_m\}_M$  by feeding mini-batch to GNN.
      // Bi-level optimization framework
      1)  $\theta_T(\{\alpha_m\}_M) \leftarrow \theta_T - lr \cdot \nabla_{\theta_T} \sum_m^M l_m$ .
      2) Update  $g(\cdot)$  through target dataset's loss function with  $\mathcal{B}_T$  and  $\theta_T(\{\alpha_m\}_M)$ .
      Update all the parameters  $\theta_T, \theta_1, \dots, \theta_M$  through  $l_T, l_1, \dots, l_M$ .
      if  $e == E$  then
        // Average affinity scores in final epoch
        for  $m \leftarrow 1, \dots, M$  do
           $\phi_m \leftarrow \phi_m + \alpha_m/I$ .
        end
      end
    end
  end
  end
  // Remove datasets according to threshold
   $\{\mathcal{D}_A\}_{M'} \leftarrow \{\mathcal{D}_{A_m} | \phi_m \geq 0.6\}$ .
   $M \leftarrow M'$ .
end
Output: Auxiliary datasets with high affinity  $\{\mathcal{D}_A\}_M$ .
```

---

33 embedding and auxiliary dataset's embedding starts with 0. The task embedding size is set as 16, and  
34 the number of the processing steps in Set2Set is set as 2. The pseudo-code is presented in Algo.1.

35 As for beam search, we apply GIN as the encoder. The beam width and search depth are both set  
36 to 3. During the search process, we train the model for 3 epochs using each candidate dataset and  
37 evaluate its performance using the validation set loss. Additionally, we consider a criterion based  
38 on the combination of the difference in fingerprint distribution and performance. Specifically, we  
39 average and normalize these two metrics to determine the criterion.

40 As for Task2vec [1], we employ GIN as the probe network and follow the official implementation<sup>3</sup>.  
41 First, we fix the encoder and train the decoder for 10 epochs. Then we apply the Monte Carlo  
42 algorithm to compute the Fisher Information Matrix.

43 As for TAG [5], we use GIN as the encoder and follow the official implementation<sup>4</sup>. TAG involves  
44 training all the datasets and computing the lookahead loss between target dataset and auxiliary  
45 datasets. The lookahead loss is accumulated over multiple epochs and used as the affinity scores.

46 As for MTDNN [10], we train all the datasets together for each target dataset and apply an additional  
47 task discriminator to classify the source of the dataset for the input instances. We train 50 epochs for  
48 GIN and 30 epochs for Graphormer, and the instances in auxiliary datasets that have a probability  
49 greater than 0.6 of being classified as the target dataset will be selected.

50 As for Gradnorm [3], we train all the datasets together and update the weights of the loss based on  
51 the gradients of the last shared GNN layer.

---

<sup>3</sup><https://github.com/awsmlabs/aws-cv-task2vec>

<sup>4</sup><https://github.com/google-research/google-research/tree/master/tag>

## 52 C Experimental Details

53 For the preliminary analysis shown in Fig.1, we conduct the study using a set of 15 molecule datasets.  
 54 Among these datasets, 11 datasets that have less than 10,000 instances in the training set are selected  
 55 as target datasets. We pair the target datasets with every other dataset and measure the relative  
 56 improvement the combination achieves. To mitigate the issue of varying dataset sizes, we upsample  
 57 or downsample the training sets of all datasets to ensure an equal number of training instances,  
 58 specifically 5,000 instances. All the reported results are based on 5 different random seeds.

59 For the dataset grouping evaluation, we train the model using the combined datasets and assess its  
 60 performance on the target datasets. We then report the model’s performance on the test set using the  
 61 best-performing model selected based on its performance on the validation set. Cross-entropy loss  
 62 is used for classification tasks and mean squared error loss is used for regression tasks. The overall  
 63 training loss is calculated as the unweighted mean of the losses for all included tasks. All the reported  
 64 results are based on 5 different random seeds.

## 65 D Overall Experimental Results

66 **Running environment.** The experiments are conducted on a single Linux server with The Intel  
 67 Xeon Gold 6240 36-Core Processor, 361G RAM, and 4 NVIDIA A100-40GB. Our method is  
 68 implemented on PyTorch 1.10.0 and Python 3.9.13.

### 69 D.1 Dataset Grouping Evaluation

70 Here we present the performance comparison over the other 5 target datasets in Table 3 and Table 4.  
 71 It can be observed that pretrained Graphormer outperforms GIN significantly, consistent with the  
 72 previous studies. In addition, MolGroup achieves the best performance in most cases, with an average  
 73 relative improvement of 3.52% and 3.10% for GIN and Graphormer. However, a notable exception  
 74 occurs with qm9 where our proposed method is unable to surpass beam search. In this instance,  
 75 MolGroup assigns low-affinity scores to each auxiliary dataset due to the significant disparity between  
 76 quantum chemistry and the other domains, as shown in Section D.3. Nonetheless, despite the limited  
 77 efficiency, beam search is capable of identifying a promising candidate by directly comparing the  
 78 performance of different groupings. It is worth noting that overall, our proposed MolGroup still  
 79 achieves better performance compared to the other baseline methods.

80 Additionally, we attribute the poor performance or even worse results of the Unweighted Average,  
 81 Gradnorm, and Pretrain-Finetune methods to the significant distribution gap between different  
 82 datasets. These methods struggle to learn a shared representation that can effectively capture the  
 83 characteristics of all the datasets.

Method	SIDER(↑)	ToxCast(↑)	ESOL(↓)	Lipo(↓)	qm9(↓)
Only-target	59.35 <sub>0.010</sub>	60.69 <sub>0.010</sub>	1.563 <sub>0.040</sub>	0.8063 <sub>0.015</sub>	0.0303 <sub>0.001</sub>
Beam search(P)	54.25 <sub>0.024</sub>	63.37 <sub>0.004</sub>	1.431 <sub>0.058</sub>	0.8073 <sub>0.011</sub>	0.0456 <sub>0.001</sub>
Beam search(P+S)	56.58 <sub>0.038</sub>	61.01 <sub>0.010</sub>	1.476 <sub>0.061</sub>	0.8147 <sub>0.014</sub>	<b>0.0258<sub>0.000</sub></b>
TAG	55.64 <sub>0.005</sub>	58.08 <sub>0.003</sub>	1.417 <sub>0.066</sub>	0.8170 <sub>0.017</sub>	0.0453 <sub>0.000</sub>
Task2vec	55.74 <sub>0.008</sub>	57.45 <sub>0.004</sub>	1.436 <sub>0.050</sub>	0.8078 <sub>0.010</sub>	0.0468 <sub>0.000</sub>
MTDNN	55.95 <sub>0.015</sub>	59.02 <sub>0.007</sub>	1.499 <sub>0.027</sub>	0.8155 <sub>0.017</sub>	0.0476 <sub>0.001</sub>
UA	56.26 <sub>0.009</sub>	59.93 <sub>0.003</sub>	1.480 <sub>0.060</sub>	0.9130 <sub>0.017</sub>	0.0494 <sub>0.000</sub>
Gradnorm	55.69 <sub>0.005</sub>	53.31 <sub>0.010</sub>	1.541 <sub>0.091</sub>	1.0755 <sub>0.001</sub>	0.0574 <sub>0.005</sub>
Pretrain-Finetune	51.43 <sub>0.009</sub>	51.06 <sub>0.005</sub>	1.480 <sub>0.099</sub>	1.0776 <sub>0.066</sub>	0.0498 <sub>0.010</sub>
MolGroup	<b>59.21<sub>0.013</sub></b>	<b>63.91<sub>0.005</sub></b>	<b>1.402<sub>0.037</sub></b>	<b>0.7996<sub>0.007</sub></b>	0.0303 <sub>0.001</sub>

Table 3: Performance comparison of GIN on target molecule datasets, with ↑ indicating higher is better and ↓ indicating lower is better.

Method	SIDER( $\uparrow$ )	ToxCast( $\uparrow$ )	ESOL( $\downarrow$ )	Lipo( $\downarrow$ )	qm9( $\downarrow$ )
Only-target	62.05 <sub>0.021</sub>	66.16 <sub>0.004</sub>	1.054 <sub>0.053</sub>	0.7432 <sub>0.032</sub>	0.0273 <sub>0.001</sub>
Beam search(P)	60.80 <sub>0.005</sub>	67.36 <sub>0.006</sub>	0.989 <sub>0.058</sub>	0.7640 <sub>0.031</sub>	0.0399 <sub>0.002</sub>
Beam search(P+S)	62.67 <sub>0.009</sub>	66.00 <sub>0.003</sub>	1.026 <sub>0.024</sub>	0.7524 <sub>0.017</sub>	<b>0.0262<sub>0.000</sub></b>
TAG	63.57 <sub>0.004</sub>	65.40 <sub>0.005</sub>	1.015 <sub>0.045</sub>	0.7507 <sub>0.014</sub>	0.0432 <sub>0.002</sub>
Task2vec	60.69 <sub>0.019</sub>	63.28 <sub>0.005</sub>	0.997 <sub>0.028</sub>	0.7562 <sub>0.015</sub>	0.0429 <sub>0.001</sub>
MTDNN	61.43 <sub>0.024</sub>	65.05 <sub>0.022</sub>	1.045 <sub>0.003</sub>	0.7716 <sub>0.032</sub>	0.0420 <sub>0.000</sub>
UW	58.24 <sub>0.014</sub>	62.71 <sub>0.020</sub>	1.120 <sub>0.084</sub>	0.7887 <sub>0.067</sub>	0.0508 <sub>0.000</sub>
Gradnorm	51.65 <sub>0.018</sub>	52.40 <sub>0.007</sub>	1.400 <sub>0.096</sub>	1.0482 <sub>0.029</sub>	0.2154 <sub>0.125</sub>
MolGroup	<b>63.75<sub>0.018</sub></b>	<b>68.68<sub>0.003</sub></b>	<b>0.978<sub>0.047</sub></b>	<b>0.7304<sub>0.018</sub></b>	0.0273 <sub>0.001</sub>

Table 4: Performance comparison using Graphormer on target molecule datasets.

## 84 D.2 Takeaway

85 **Grouping more high-affinity datasets improves performance.** Previous studies on task grouping  
86 have assumed that low-order relationships can be an effective indicator of high-order ones. It  
87 suggests that combining multiple source datasets, which can benefit the target dataset, leads to better  
88 performance when learned together. Our experimental results also confirm this phenomenon. We take  
89 eight datasets as examples and add them one by one with the top three datasets having the highest  
90 affinity score as measured by MolGroup. Results are presented in Table 5, where  $\text{Top}\{a, b, \dots\}$   
91 denotes the combination of auxiliary datasets with top-performing ones. We can find that combining  
92 more auxiliary datasets leads to better performance in most cases. Besides, training with high-affinity  
93 datasets can significantly reduce the variant of FreeSolv (1.579  $\rightarrow$  0.279), indicating a more robust  
94 representation learned from the auxiliary datasets.

Combinations	ClinTox	Tox21	FreeSolv	BBBP	BACE	ToxCast	ESOL	Lipo
Only Target	56.45 <sub>0.023</sub>	74.23 <sub>0.005</sub>	3.842 <sub>1.579</sub>	66.62 <sub>0.028</sub>	75.02 <sub>0.026</sub>	60.69 <sub>0.010</sub>	1.563 <sub>0.040</sub>	0.8063 <sub>0.015</sub>
+ Top{1}	57.77 <sub>0.028</sub>	74.81 <sub>0.006</sub>	3.563 <sub>0.989</sub>	67.36 <sub>0.023</sub>	71.27 <sub>0.045</sub>	62.66 <sub>0.002</sub>	1.502 <sub>0.043</sub>	0.8096 <sub>0.014</sub>
+ Top{1,2}	57.48 <sub>0.032</sub>	75.22 <sub>0.003</sub>	3.462 <sub>0.970</sub>	<b>68.64<sub>0.012</sub></b>	71.13 <sub>0.019</sub>	63.18 <sub>0.006</sub>	1.524 <sub>0.077</sub>	0.8078 <sub>0.011</sub>
+Top{1,2,3}	<b>59.77<sub>0.027</sub></b>	<b>75.66<sub>0.004</sub></b>	<b>3.116<sub>0.279</sub></b>	68.36 <sub>0.016</sub>	<b>77.33<sub>0.015</sub></b>	<b>63.91<sub>0.005</sub></b>	<b>1.402<sub>0.010</sub></b>	<b>0.7996<sub>0.005</sub></b>

Table 5: Performance of top-k combinations.

95 **PCBA is an effective booster.** One interesting finding is that dataset PCBA [6] can boost the  
96 performance of most of the small molecule datasets, as shown in Table 6. This dataset offers both  
97 a diverse range of chemical compounds with unique scaffold structures, comprising over 350,000  
98 training instances, and an extensive collection of 128 bioassay annotations that represent a broad  
99 range of biological activities, making it a potent booster for small molecule property prediction tasks  
100 that can benefit from both structure and task.

	BBBP( $\uparrow$ )	ClinTox( $\uparrow$ )	ToxCast( $\uparrow$ )	Tox21( $\uparrow$ )	ESOL( $\downarrow$ )	FreeSolv( $\downarrow$ )	Lipo( $\downarrow$ )
Only-target	66.62 <sub>0.028</sub>	56.45 <sub>0.023</sub>	60.69 <sub>0.010</sub>	74.23 <sub>0.005</sub>	1.563 <sub>0.040</sub>	3.842 <sub>1.579</sub>	0.8063 <sub>0.015</sub>
+PCBA	67.11 <sub>0.023</sub>	57.77 <sub>0.028</sub>	62.05 <sub>0.007</sub>	74.81 <sub>0.006</sub>	1.463 <sub>0.020</sub>	3.563 <sub>0.989</sub>	0.8021 <sub>0.009</sub>

Table 6: Cases whose performance is improved by PCBA.

## 101 D.3 Case Study

102 To give an intuition of the learning process of MolGroup, we show the learning curves of the affinity  
103 scores in different filtering rounds in Fig. 1. In each round, auxiliary datasets with affinity scores  
104 below 0.6 are removed, and if none of the datasets fall below this threshold, the dataset with the  
105 lowest affinity score is filtered out. It can be observed that a significant number of auxiliary datasets  
106 are removed in the first or second round. Furthermore, the learning curves tend to converge after  
107 6 epochs in most cases. The dataset PCBA remains in the final round in most cases, indicating its  
108 general benefit to the other target datasets. We also notice that the majority of auxiliary datasets are

assigned high-affinity scores for FreeSolv, as demonstrated in the preliminary experiment in Section 1, suggesting that all auxiliary datasets contribute positively to its performance.

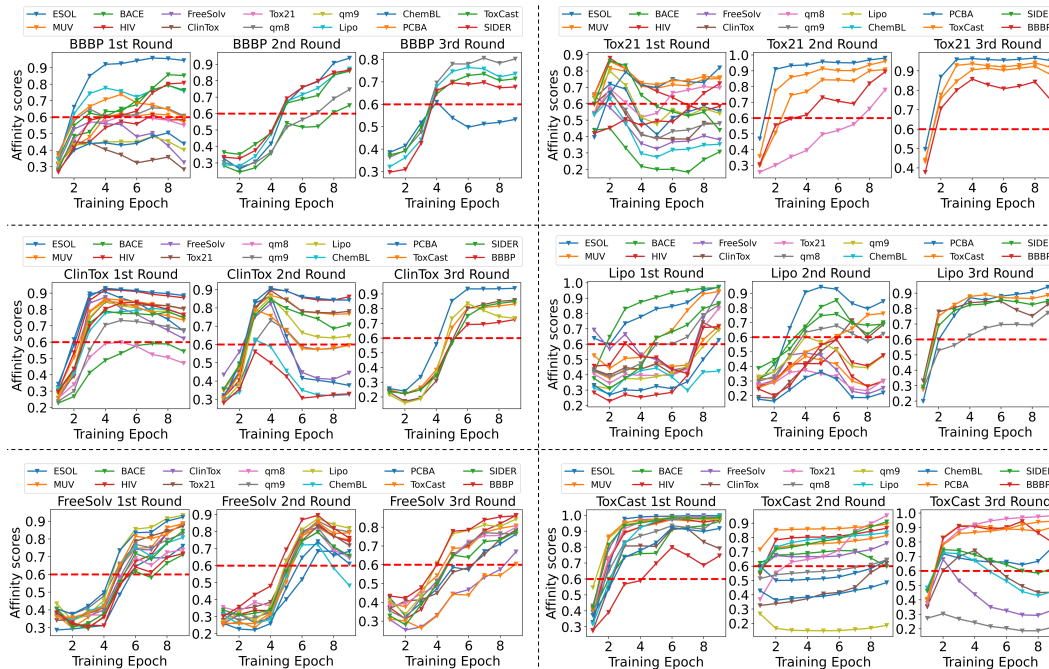


Figure 1: Learning curves of affinity scores where red dashed line represents threshold.

Additionally, we plot the learning curves for qm8 and qm9, which don't have suggested auxiliary datasets. As shown in Fig. 2, all the auxiliary datasets are assigned affinity scores lower than the threshold for both qm8 and qm9, resulting in the removal of all datasets after the first round. We attribute this to the large discrepancy in the structure and task between quantum chemistry and the other domains such as medication. Molecules in quantum chemistry have diverse structures, including both natural and hypothetical compounds. They are studied to explore the behaviors of various functional groups. On the other hand, molecules in medication are more focused on structure. Their study revolves around predicting and optimizing molecular properties relevant to drug design.

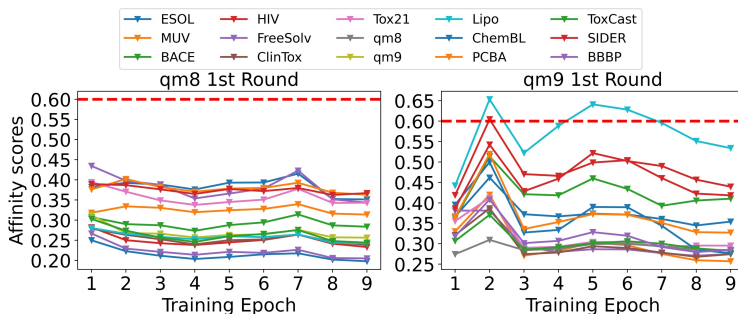


Figure 2: Learning curves for qm8 and qm9.

## E Parameter Analysis

### E.1 Analysis on balance coefficient $\lambda$

To analyze the impact of the balance coefficient  $\lambda$ , we vary  $\lambda$  in the range of  $\{0.9, 0.7, 0.5, 0.3\}$  and pick BBBP, Tox21, and Lipo as examples. The learning curves of the affinity scores corresponding to different values of  $\lambda$  are plotted in Fig. 3. From the results, we observe that decreasing  $\lambda$  led to lower affinity scores for the auxiliary datasets. However, these scores fail to effectively discriminate the

125 affinity of individual auxiliary datasets. One reason for this is the instability introduced by parameter  
 126 initialization and the per-step level computation of structure affinity scores. As a result, lower values  
 127 of  $\lambda$  cause the MolGroup to assign lower affinity scores to stabilize the training of the target dataset.  
 128 In light of this, we suggest setting  $\lambda$  to a high value, i.e., 0.9.

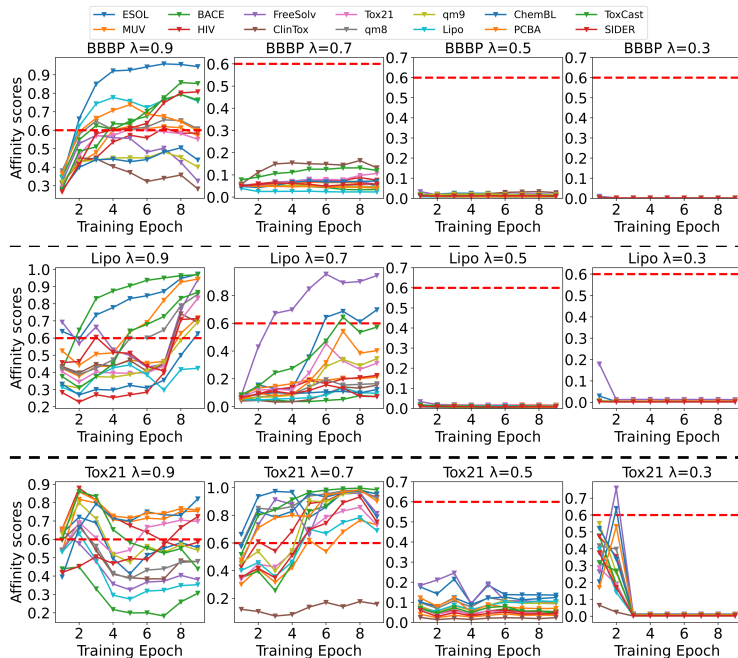


Figure 3: Learning curves of affinity scores with different  $\lambda$ .

## 129 E.2 Analysis on number of filtering rounds $R$

130 As shown in Fig. 1, the auxiliary datasets with the highest affinity scores change as we progress  
 131 through the different filtering rounds. To investigate the impact of the filtering rounds, we test the  
 132 performance of the top-3 auxiliary datasets selected in each rounds. The comparison results are  
 133 shown in Table 7 and the selected datasets are illustrated in Fig. 4. We observe that the auxiliary  
 134 datasets selected in the 3rd round exhibit the best performance. This is attributed to the filtering  
 135 process in the 1st and 2nd rounds, which removes negative datasets and helps alleviate interference.  
 136 As a result, MolGroup can estimate the affinity of each auxiliary dataset more accurately, leading to  
 137 improved performance on the target dataset.

	BBBP	Tox21	ClinTox	Lipo	FreeSolv	ToxCast
$R = 1$	67.94 <sub>0.018</sub>	75.66 <sub>0.004</sub>	57.66 <sub>0.027</sub>	0.8013 <sub>0.013</sub>	3.2288 <sub>0.543</sub>	58.82 <sub>0.007</sub>
$R = 2$	67.99 <sub>0.018</sub>	75.66 <sub>0.004</sub>	57.66 <sub>0.027</sub>	0.8020 <sub>0.031</sub>	3.2280 <sub>0.627</sub>	63.47 <sub>0.003</sub>
$R = 3$	<b>68.36</b> <sub>0.016</sub>	<b>75.66</b> <sub>0.004</sub>	<b>59.77</b> <sub>0.027</sub>	<b>0.7996</b> <sub>0.007</sub>	<b>3.116</b> <sub>0.279</sub>	<b>63.91</b> <sub>0.005</sub>

Table 7: Performance with different number of filtering rounds  $R$ .

## 138 F Broader impact

139 **Impact on machine learning research.** We propose a novel strategy that combines the routing  
 140 mechanism with the meta gradient to quantify the impact of one dataset on another. Previously, the  
 141 routing mechanism was used to increase the model capacity. Our proposed framework can inspire  
 142 various extensive applications in machine learning, including neural architecture searching (NAS)  
 143 and data-centric AI. Specifically, our framework enables the network to determine the optimal routing  
 144 path through the meta gradient. This empowers the network to control and modify the layerwise  
 145 architecture, leading to improved performance. Besides, it has the potential to enhance data-centric



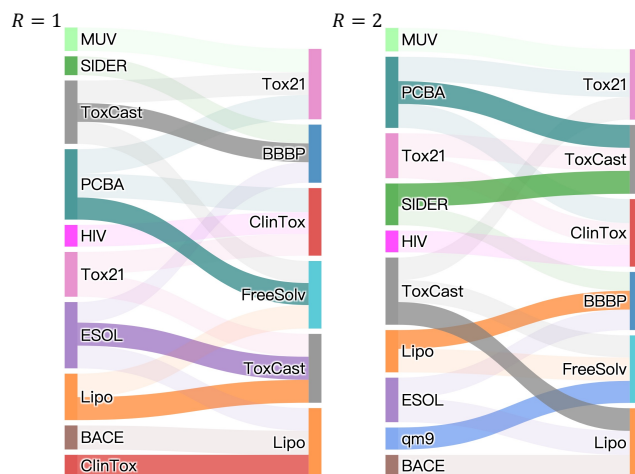


Figure 4: Auxiliary datasets with top-3 affinity scores with different  $R$  where we highlight the different edges in these two rounds.

AI approaches by providing a tool to analyze and understand the relationship between different datasets or sub-datasets.

**Impact on biological research.** We investigate the relationship between molecule datasets, and, focusing on both the structural and task dimensions. our findings shed light on how different molecule datasets impact each other. The analytical approach we employ can be extended to other biological data, such as protein and RNA. Furthermore, our method has the potential to be utilized for data instance filtering in the biological domain. This is particularly important given that biological data often contains various types of noise. By filtering out data instances with negative effects on downstream tasks, we can improve the quality and reliability of biological data used in research.

**Impact on the society.** Our study has significant implications for society, particularly in the field of biomedicine and drug discovery. By understanding the relationship between molecule datasets and their impact on each other, we can gain insights into how different molecules interact and influence biological processes. It can aid in the development of more effective drugs and therapies by identifying molecules that have a positive impact on specific biological targets or diseases.

## References

- [1] Alessandro Achille, Michael Lam, Rahul Tewari, Avinash Ravichandran, Subhansu Maji, Charless C Fowlkes, Stefano Soatto, and Pietro Perona. Task2vec: Task embedding for meta-learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6430–6439, 2019.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, and Andrew Rabinovich. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *International conference on machine learning*, pages 794–803. PMLR, 2018.
- [4] Eli Fernández-de Gortari, César R García-Jacas, Karina Martínez-Mayorga, and José L Medina-Franco. Database fingerprint (dfp): an approach to represent molecular databases. *Journal of cheminformatics*, 9(1):1–9, 2017.
- [5] Chris Fifty, Ehsan Amid, Zhe Zhao, Tianhe Yu, Rohan Anil, and Chelsea Finn. Efficiently identifying task groupings for multi-task learning. *Advances in Neural Information Processing Systems*, 34:27503–27516, 2021.
- [6] Weihua Hu, Matthias Fey, Hongyu Ren, Maho Nakata, Yuxiao Dong, and Jure Leskovec. Ogb-lsc: A large-scale challenge for machine learning on graphs. *arXiv preprint arXiv:2103.09430*, 2021.
- [7] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems*, 33:22118–22133, 2020.
- [8] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015.
- [9] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [10] Po-Nien Kung, Sheng-Siang Yin, Yi-Cheng Chen, Tse-Hsuan Yang, and Yun-Nung Chen. Efficient multi-task auxiliary learning: selecting auxiliary data by feature similarity. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 416–428, 2021.
- [11] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [12] David Mendez, Anna Gaulton, A Patrícia Bento, Jon Chambers, Marleen De Veij, Eloy Félix, María Paula Magariños, Juan F Mosquera, Prudence Mutowo, Michał Nowotka, et al. ChEMBL: towards direct deposition of bioassay data. *Nucleic acids research*, 47(D1):D930–D940, 2019.
- [13] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [14] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):513–530, 2018.
- [15] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.
- [16] Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? *Advances in Neural Information Processing Systems*, 34:28877–28888, 2021.