



Improving Out-of-Distribution Generalization in Graphs via Hierarchical Semantic Environments

 Biomedical domain

Yinhua Piao, Sangseon Lee, Yijingxiu Lu, Sun Kim.

Bio & Health Informatics Lab

Department of Computer Science and Engineering

Seoul National University, Seoul, Korea



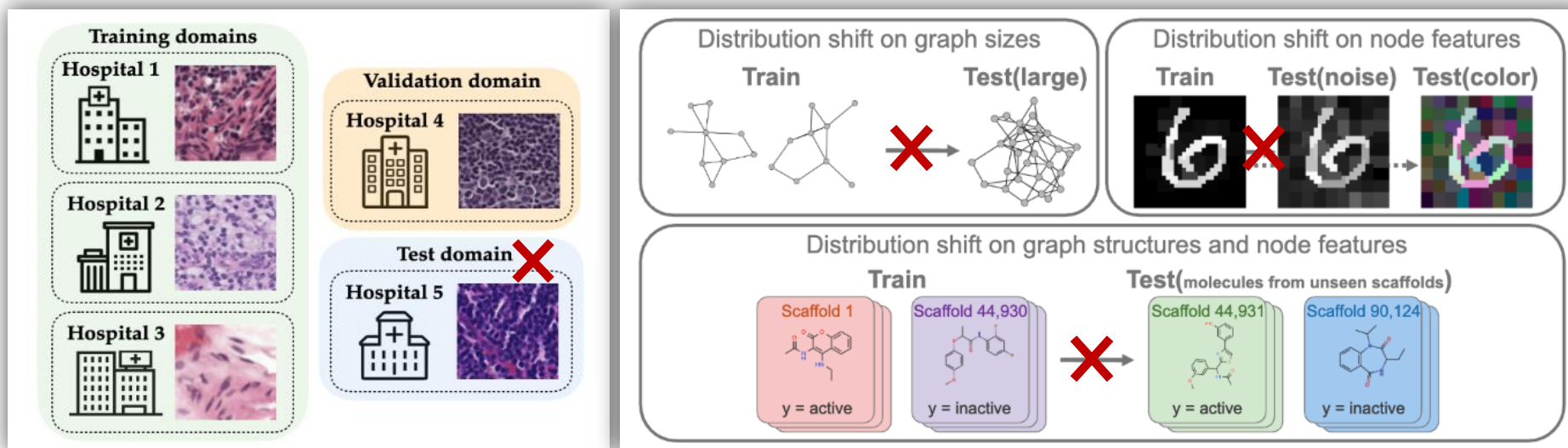
Out-of-distribution Generalization

Out-of-distribution Generalization

- Models learned with **Empirical Risk Minimization** often **fail** to generalize to OOD data.

$$\min_f \mathbb{E}_{(x,y) \sim P_{\text{train}}(x,y)} [l(f(x), y)], \quad P_{\text{test}}(x,y) \neq P_{\text{train}}(x,y)$$

- Applications: Molecule prediction, pandemic prediction, medical detection for COVID-19.



(A robe, et al., 2021,
Li, Haoyang, et al., 2022,
Arjovsky et al., 2022)

OOD Generalization in Euclidean and Non-Euclidean Domains

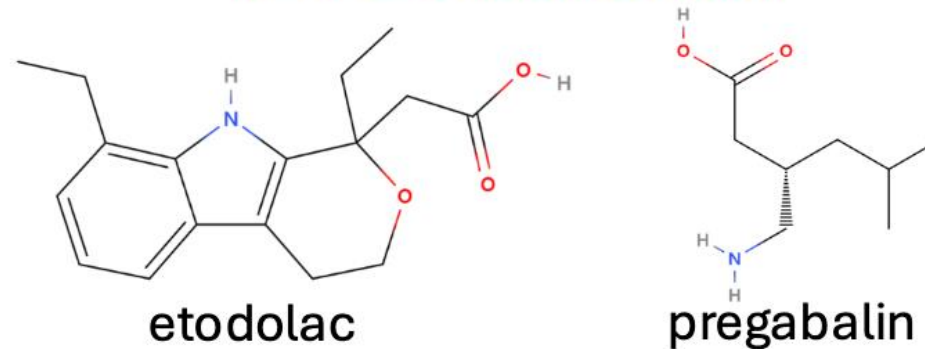
Euclidean-based OOD: feature distribution shifts

Non-euclidean-based OOD: feature & structural shifts simultaneously.

Euclidean domain



Non-Euclidean domain



Euclidean
domain

Feature

Structure

Non-Euclidean
domain

In-distribution

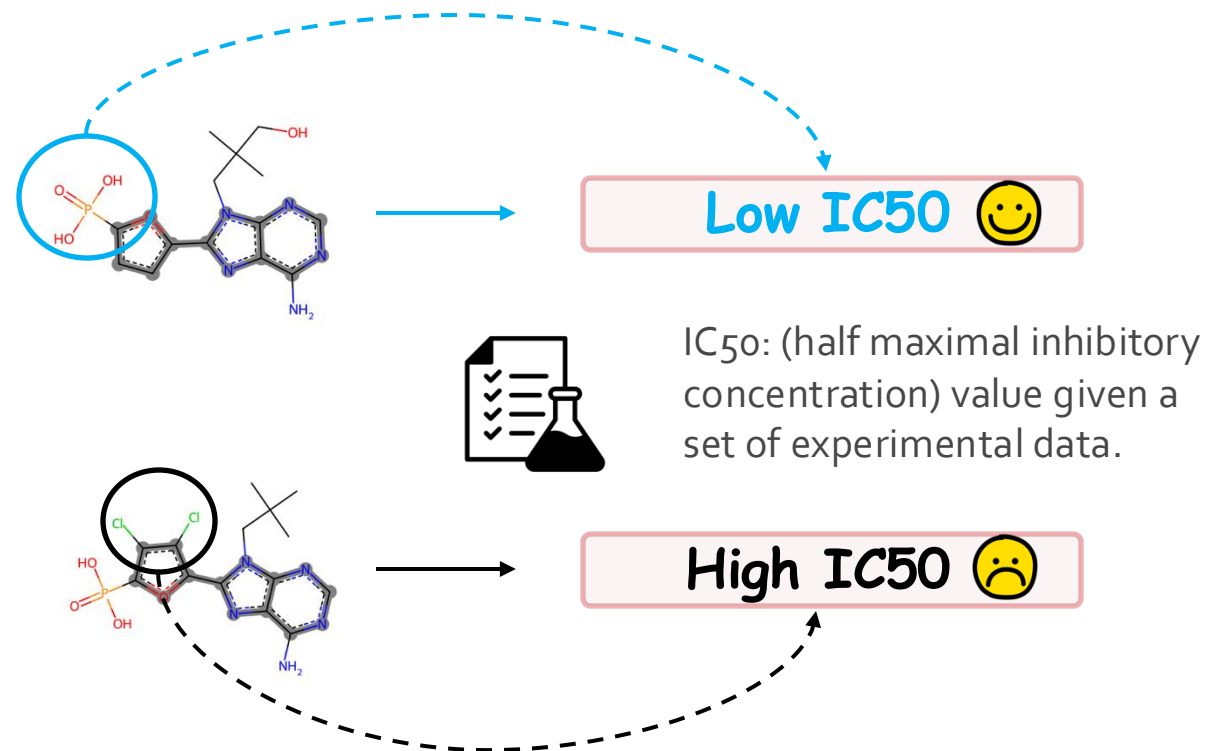
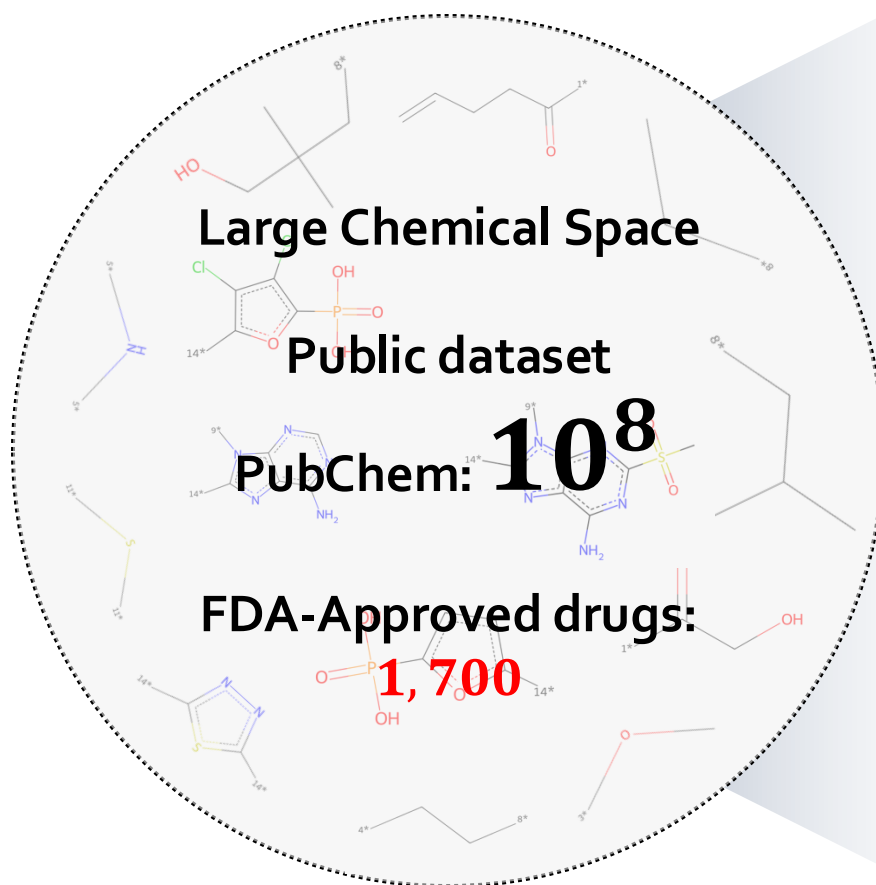
shifts

Out-of-distribution



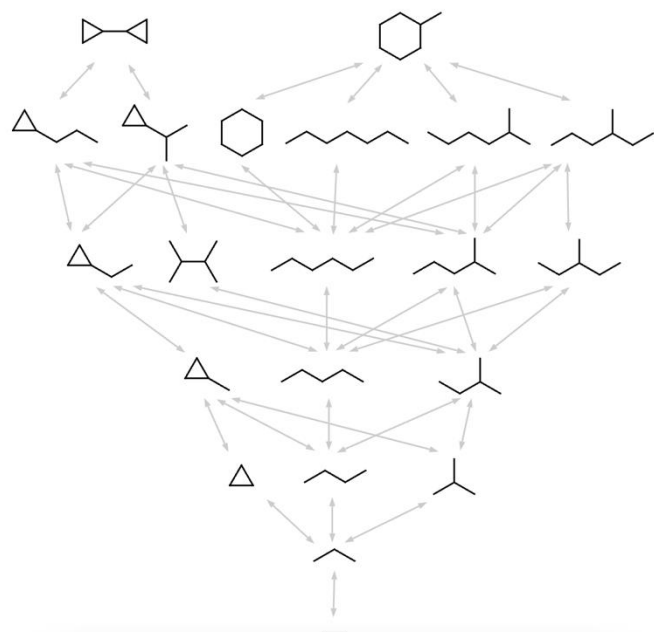
Graph OOD Generalization in Chemical Domain

Molecules and Molecular Property Prediction



Challenge of Subgraph Mining

➤ Subgraph mining



Subgraph is unknown !

For a molecule with n bonds, number of possible subgraphs:

2^n combinations

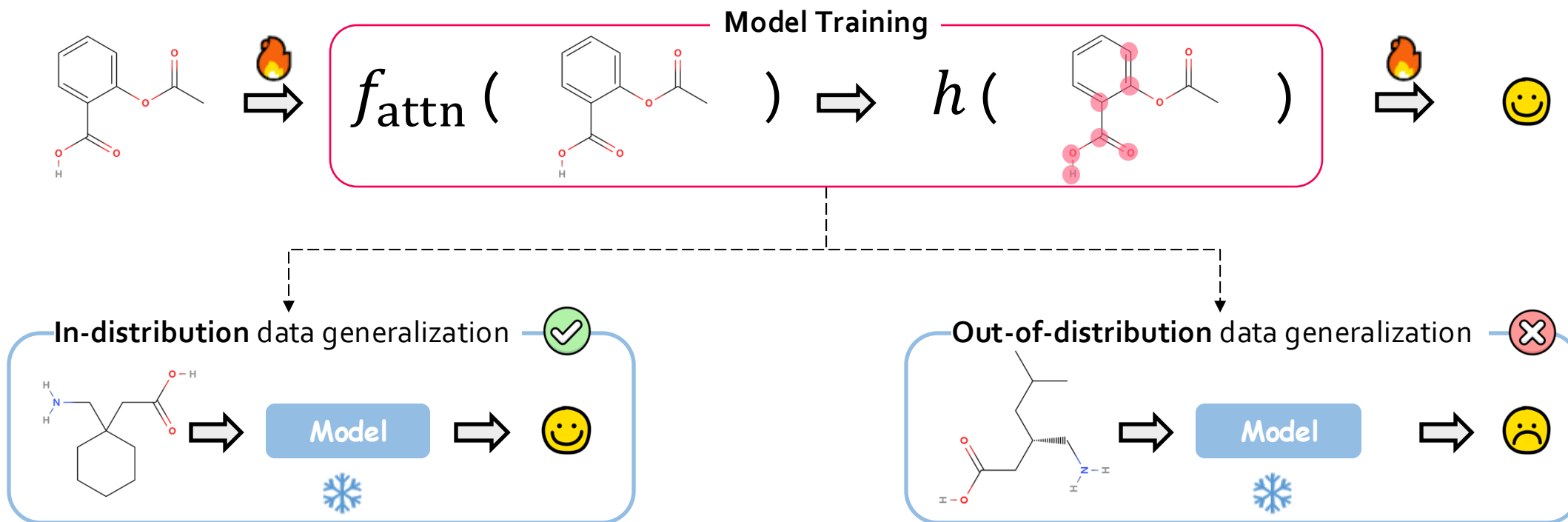
(exponential growth)



Huge chemical subgraph search space

*SmallWorld: Efficient Maximum Common Subgraph Searching of Large Chemical Databases.
ACS National Meeting, Philadelphia, USA 22nd August 2012*

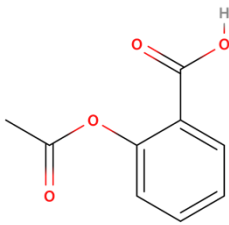
Subgraph modelling is important in graph OOD generalization



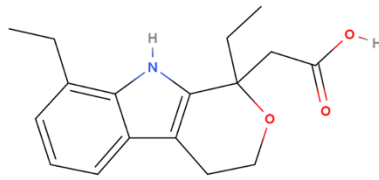
Poor generalization on out-of-distribution data.

An example

Train data
(Ring structured)



Aspirin

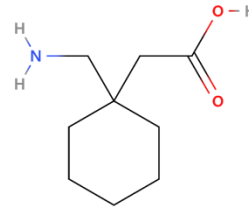


Etodolac

Model 

Relieve pain

In-distribution Test data
(Ring-structured)

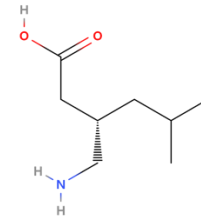


Gabapentin

Model 

Relieve pain ?

Out-of-distribution Test data
(non-Ring structured)

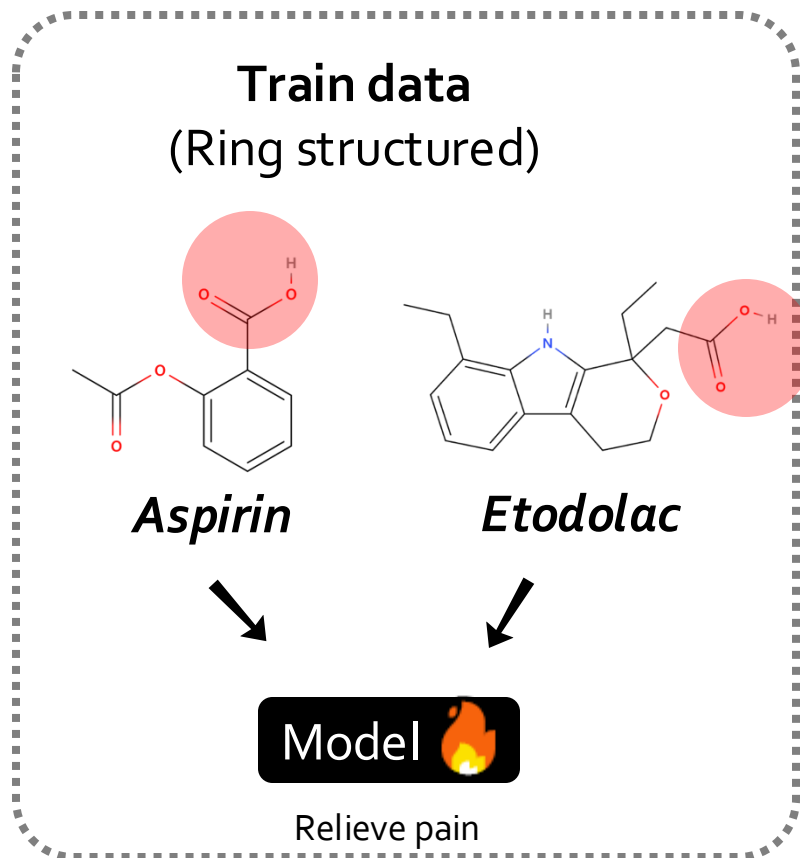


Pregabalin

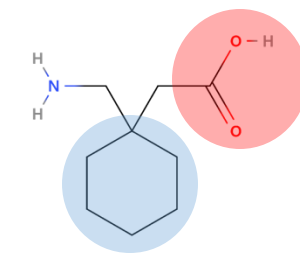
Model 

Relieve pain ?

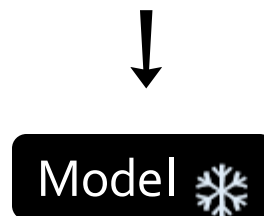
An example



In-distribution Test data
(Ring-structured)

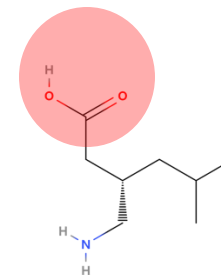


Gabapentin

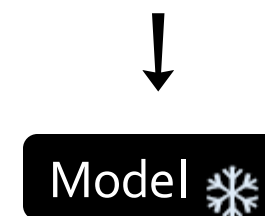


Relieve pain ✓

Out-of-distribution Test data
(non-Ring structured)



Pregabalin



Others ✗

NO RING?

✗ More challenging on unseen and highly different molecules (size, scaffolds, etc)

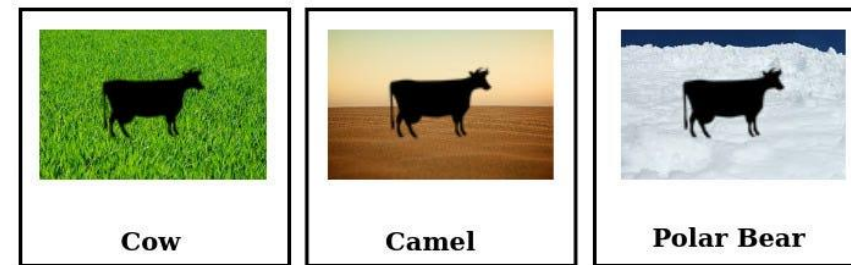
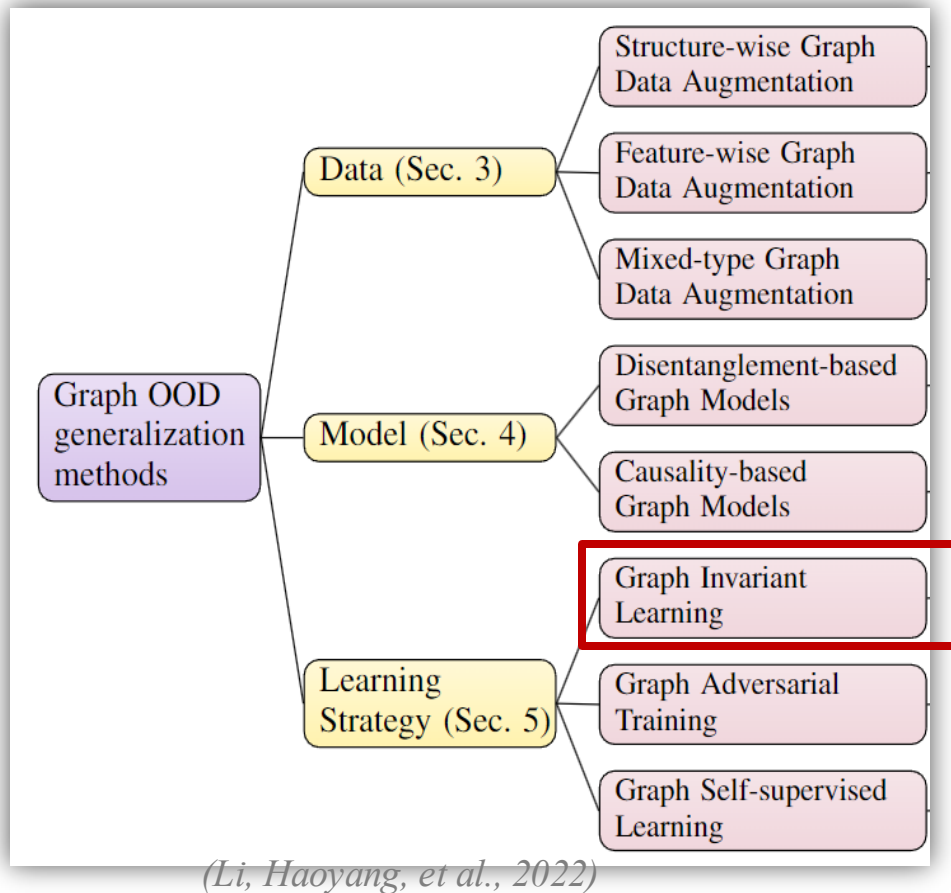
? How to learn an effective subgraph for the unknown molecule's property?



Related works on Graph OOD Generalization

Graph Invariant Learning

Invariant Learning treats the cause of distribution shifts between testing and training data as a potential unknown environmental variable $\mathbf{e} \in \mathcal{E}$ (*hospital, size, color, scaffold, etc*).



Neural Network Predictions

Empirical Risk Minimization

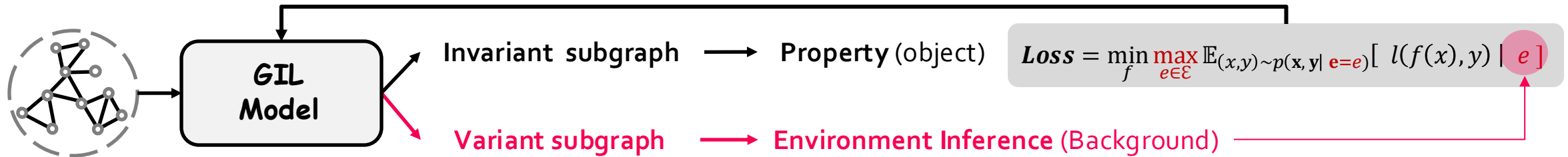
$$\min_f \mathbb{E}_{(x,y) \sim p(x,y)} [l(f(x), y)]$$

Invariant Risk Minimization

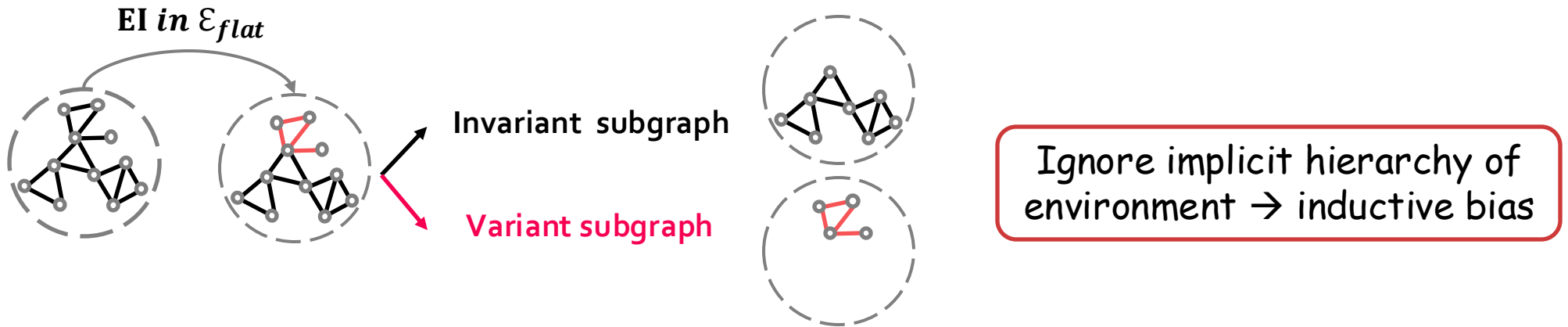
$$\min_f \max_{\mathbf{e} \in \mathcal{E}} \mathbb{E}_{(x,y) \sim p(x,y | \mathbf{e}=\mathbf{e})} [l(f(x), y) | \mathbf{e}]$$

Graph Invariant Learning – Existing works

- Graph Invariant Learning (GIL) improves Model generalization on OOD data.

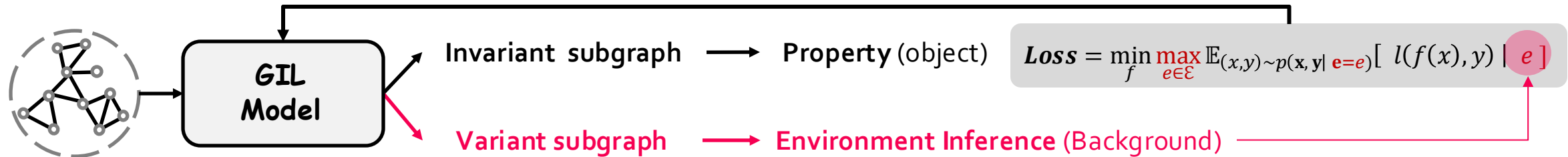


Existing GIL Method: Flat environment inference.

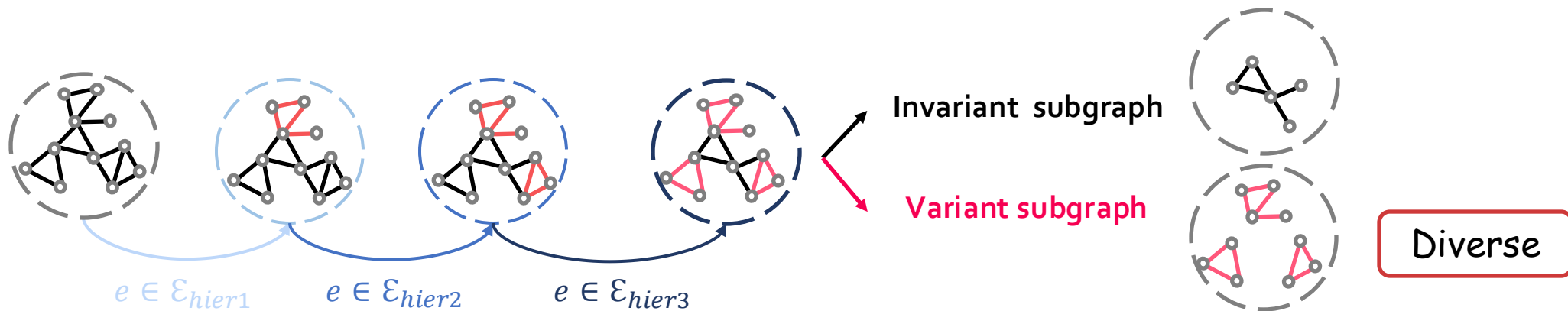


Graph Invariant Learning – Our works

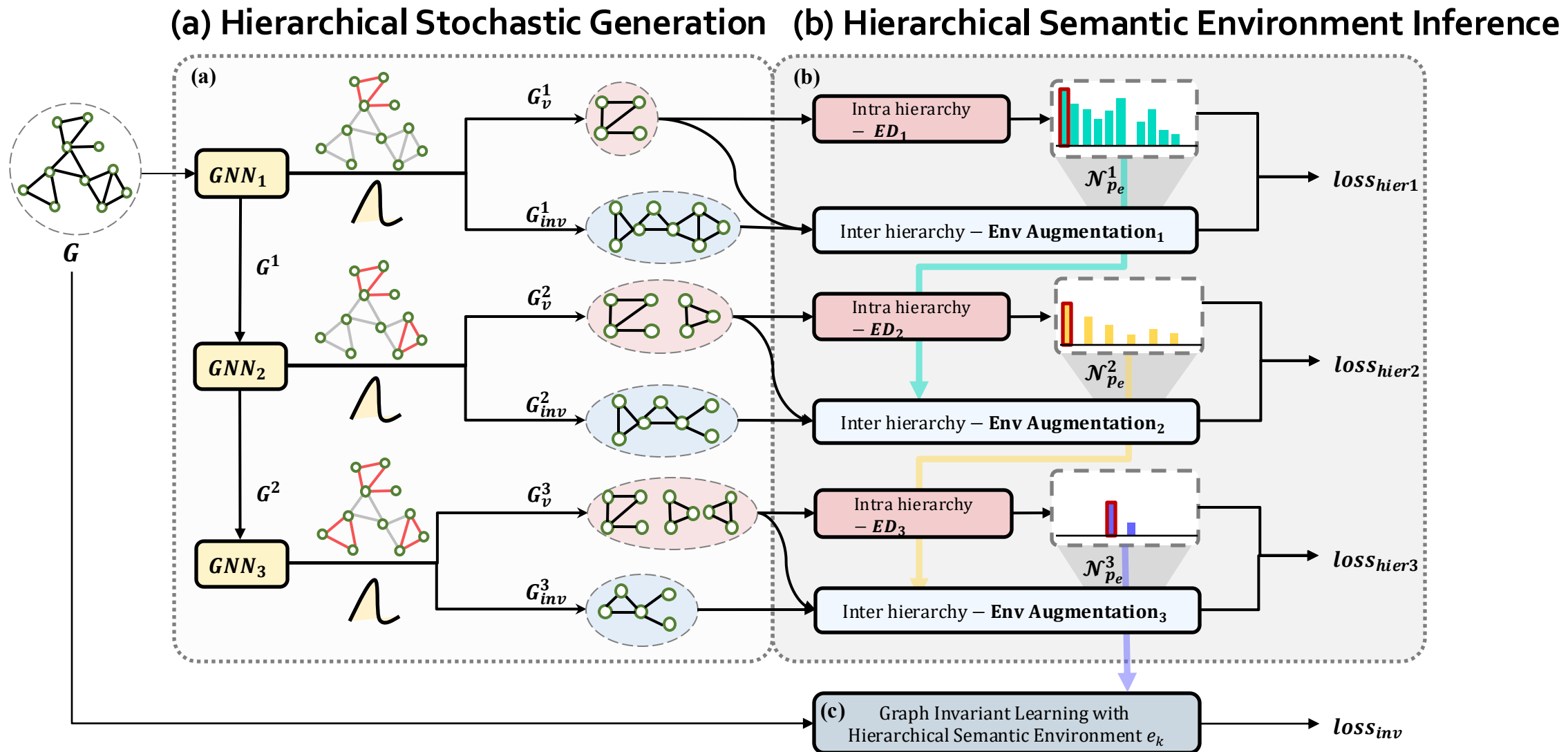
- Graph Invariant Learning (GIL) improves Model generalization on OOD data.



Our Method: Hierarchical environment inference.



Our Approach

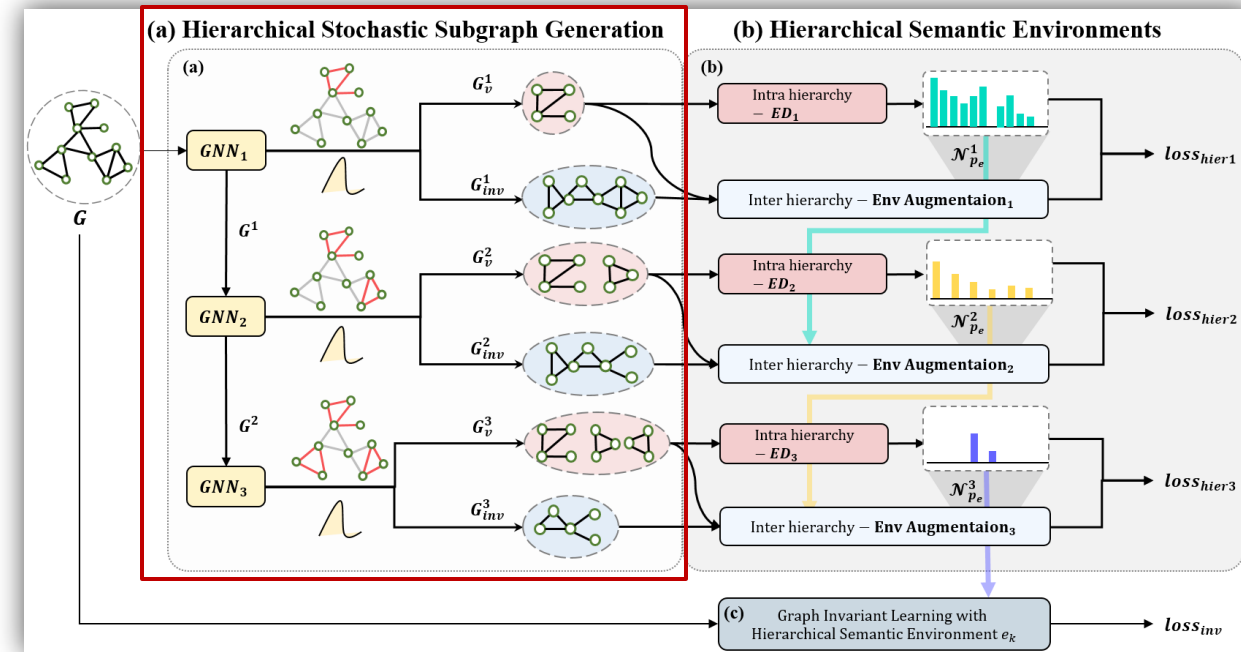


Our Approach

(a) Hierarchical Stochastic Generation

- Input: $G = (V, A)$, $h = \text{GNN}(V, A)$
 - Probability distribution of edge:
$$s_{ij} = S(h_{ij}) = \sigma(\text{MLP}([h_i, h_j]))$$
- Stochastic Edge Selection via Gumbel-softmax:
$$\mathbf{p}_{ij} \in \{0,1\} \sim \text{Bern}(s_{ij})$$
- Hierarchical Stochastic Subgraph Generation
$$N_{ij}^k = N_{ij}^{k-1} + 1\{N_{ij}^{k-1} = 0 \text{ and } \mathbf{p}_{ij}^k > T\},$$

$$A_v^k \leftarrow A \odot N^k, \quad A_{inv}^k \leftarrow A - A_v^k$$
- Output: $\{G_v^k: \text{variant subgraph}, G_{inv}^k: \text{invariant subgraph}\}$



Our Approach

(b) Hierarchical Semantic Environments

- Intra-Hierarchy Environment Diversification.

$$L_{ED} = -\frac{1}{K} \sum_k \sum_{e_k} \max_{e_k} \log(P(\mathbf{e}_k | f^e(G_{v_k}, y)))$$

Variant subgraph

- Inter-Hierarchy Environment Augmentation

$$L_{EnvCon} = \mathbb{E}_{x \sim p_d} \left[\frac{1}{K} \sum_k L_{InfoNCE}(z, N_{p_e}^k(z), \tau) \right],$$

$$L_{LabelCon} = \mathbb{E}_{x \sim p_d} \left[\frac{1}{K} \sum_k L_{InfoNCE}(z, N_{p_y}(z), \tau) \right], \text{ where}$$

$$N_{p_e}^k(z) = N_{p_e}^{k-1}(z) \cup \{z_i \mid e_i^k = e_z^k, z_i \in N(z)\},$$

$$N_{p_y}(z) = \{z_i \mid y_i = y_z, z_i \in N(z)\}$$

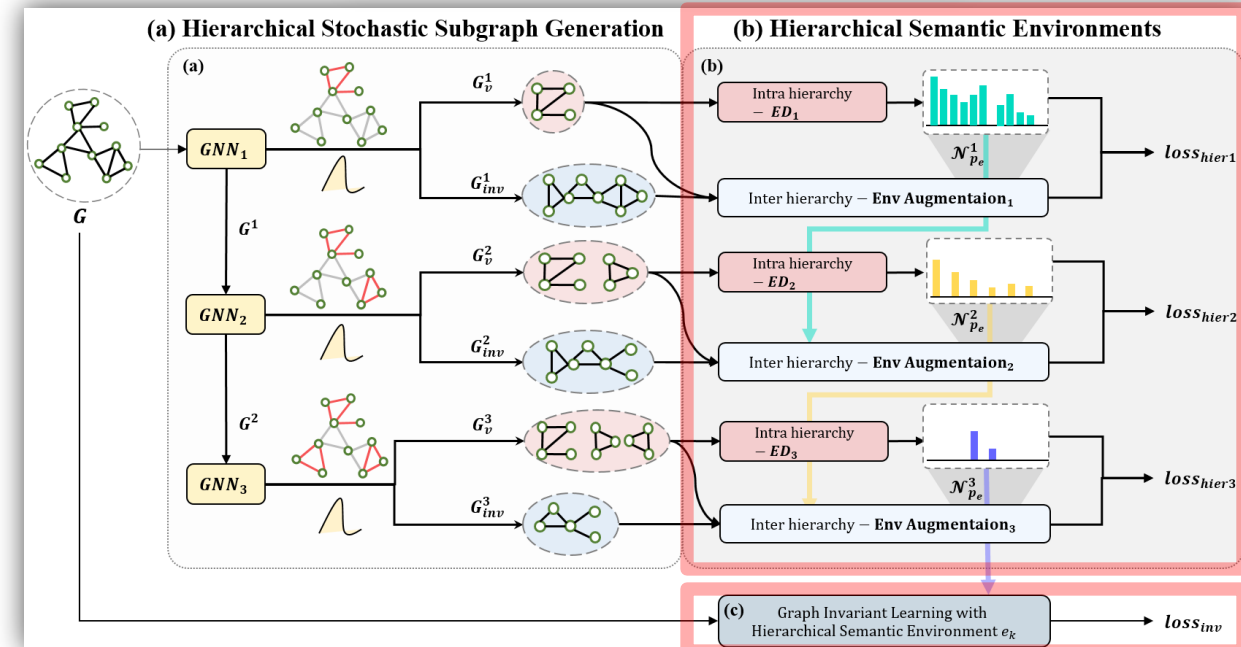
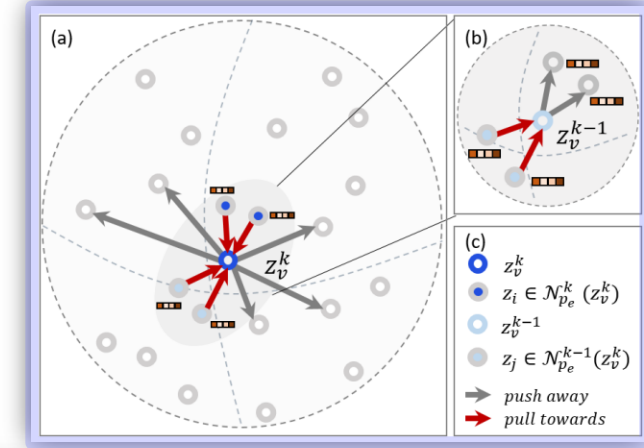
(c) Graph invariant learning with HSE

- Hierarchical Semantic Environment Learning

$$L_{HEI} = L_{ED} + \alpha \cdot L_{EnvConv} + \beta \cdot L_{LabelCon}$$

- Robust GIL with HSE at k-th hierarchy.

$$\min_f L_{cls}^{e_k}(f) + Var(L_{cls}^{e_k}(f)) \text{ s.t. } \mathbf{e}_k = \arg \min_{e_k} L_{HEI}$$

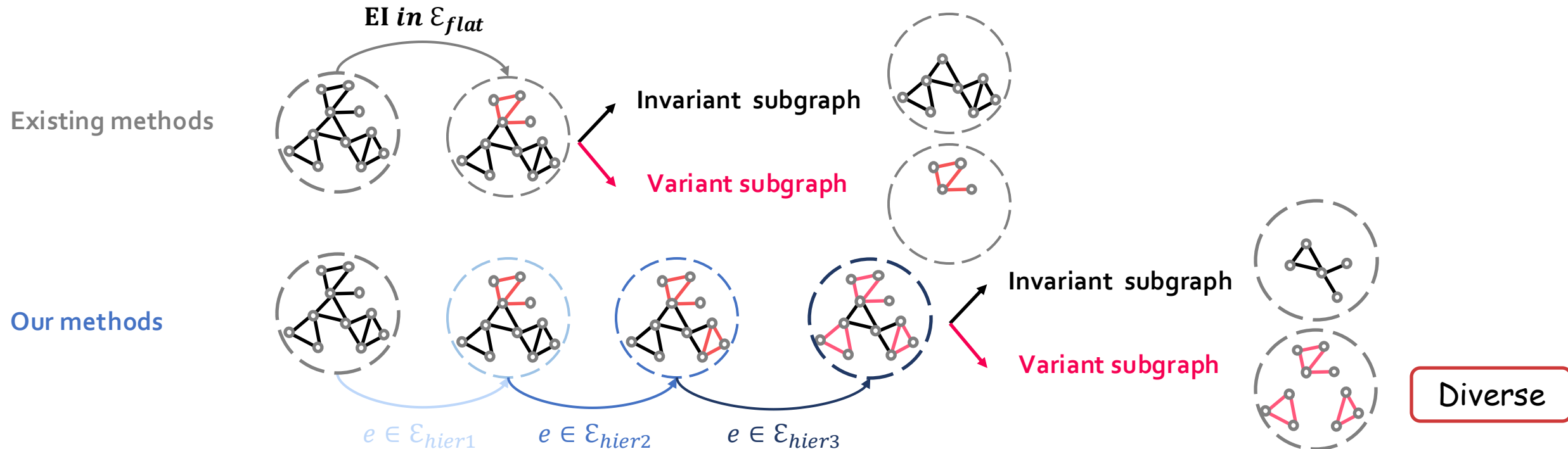


Flat Env Inference vs Hierarchical Env Inference

Two extreme cases in Flat Env Inference approaches:

1. Too Few Environments \rightarrow embedding space collapse.
2. Too many individual Environments \rightarrow ignore dependency between environments.

-> can also view as an adaptive parameter search problem.





Experiments and Results

Results

- Compared to ERM, Euclidean-based invariant learning show degraded performance.
- Graph-based OOD methods exhibit better performance.
- Our method outperforms all baselines, especially when comes to complex dataset (*-SCA).

Table 1. Statistics of the datasets used in experiments. The number of nodes and edges are averaged numbers among all datasets.

| DATASETS | #TRAINING | #VALIDATION | #TESTING | #LABELS | #ENVS | #NODES | #METRICS |
|------------|-----------|-------------|----------|---------|-------|--------|----------|
| CMINST-SP | 40,000 | 5,000 | 15,000 | 2 | N.A. | 56.90 | ACC |
| GRAPH-SST5 | 6,090 | 1,186 | 2,240 | 5 | N.A. | 19.85 | ACC |
| IC50-ASSAY | 34,179 | 19,028 | 19,032 | 2 | 311 | 34.58 | ROC-AUC |
| IC50-SCA | 21,519 | 19,041 | 19,048 | 2 | 6,881 | 39.38 | ROC-AUC |
| IC50-SIZE | 36,597 | 17,660 | 16,415 | 2 | 190 | 37.99 | ROC-AUC |
| EC50-ASSAY | 4,540 | 2,572 | 2,490 | 2 | 47 | 39.81 | ROC-AUC |
| EC50-SCA | 2,570 | 2,532 | 2,533 | 2 | 850 | 56.84 | ROC-AUC |
| EC50-SIZE | 4,684 | 2,313 | 2,398 | 2 | 167 | 48.40 | ROC-AUC |

Table 2. Test ROC-AUC of various models on DrugOOD benchmark datasets. The mean \pm standard deviation of all models is reported as an average of 5 executions of each model. The best methods are highlighted in bold and the second best methods are underlined.

| METHODS | IC50-ASSAY | IC50-SCA | IC50-SIZE | EC50-ASSAY | EC50-SCA | EC50-SIZE |
|--------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| ERM[51] | 71.79 \pm 0.27 | 68.85 \pm 0.62 | 66.70 \pm 1.08 | 76.42 \pm 1.59 | 64.56 \pm 1.25 | 62.79 \pm 1.15 |
| IRM[23] | 72.12 \pm 0.49 | 68.69 \pm 0.65 | 66.54 \pm 0.42 | 76.51 \pm 1.89 | 64.82 \pm 0.55 | 63.23 \pm 0.56 |
| V-REX[25] | 72.05 \pm 1.25 | 68.92 \pm 0.98 | 66.33 \pm 0.74 | 76.73 \pm 2.26 | 62.83 \pm 1.20 | 59.27 \pm 1.65 |
| EIIL[26] | 72.60 \pm 0.47 | 68.45 \pm 0.53 | 66.38 \pm 0.66 | 76.96 \pm 0.25 | 64.95 \pm 1.12 | 62.65 \pm 1.88 |
| IB-IRM[24] | 72.50 \pm 0.49 | 68.50 \pm 0.40 | 66.64 \pm 0.28 | 76.72 \pm 0.98 | 64.43 \pm 0.85 | 64.10 \pm 0.61 |
| GREASE [36] | 72.77 \pm 1.25 | 68.33 \pm 0.32 | 66.16 \pm 0.46 | 72.44 \pm 2.55 | <u>67.98\pm1.00</u> | 63.93 \pm 3.01 |
| CIGAV1 [33] | 72.71 \pm 0.52 | 69.04 \pm 0.86 | 67.24 \pm 0.88 | 78.46 \pm 0.45 | 66.05 \pm 1.29 | <u>66.01\pm0.84</u> |
| CIGAV2 [33] | <u>73.17\pm0.39</u> | <u>69.70\pm0.27</u> | <u>67.78\pm0.76</u> | - | - | - |
| MOLEOOD [31] | 71.38 \pm 0.68 | 68.02 \pm 0.55 | 66.51 \pm 0.55 | 73.25 \pm 1.24 | 66.69 \pm 0.34 | 65.09 \pm 0.90 |
| GALA [34] | - | - | - | <u>79.24\pm1.36</u> | 66.00 \pm 1.86 | <u>66.01\pm0.84</u> |
| OURS | 74.01\pm0.11 | 70.72\pm0.30 | 68.64\pm0.23 | 80.82\pm0.21 | 69.73\pm0.21 | 66.87\pm0.38 |

Results

➤ Discussion on Diversity of Inferred Environments ➤ Case study of Hierarchical Semantic Environment.

(a) Random sampling methods.

(b) Flat environment inference methods

(c) Our hierarchical environment inference methods

(d) K-S test to calculate the diversity of three methods.

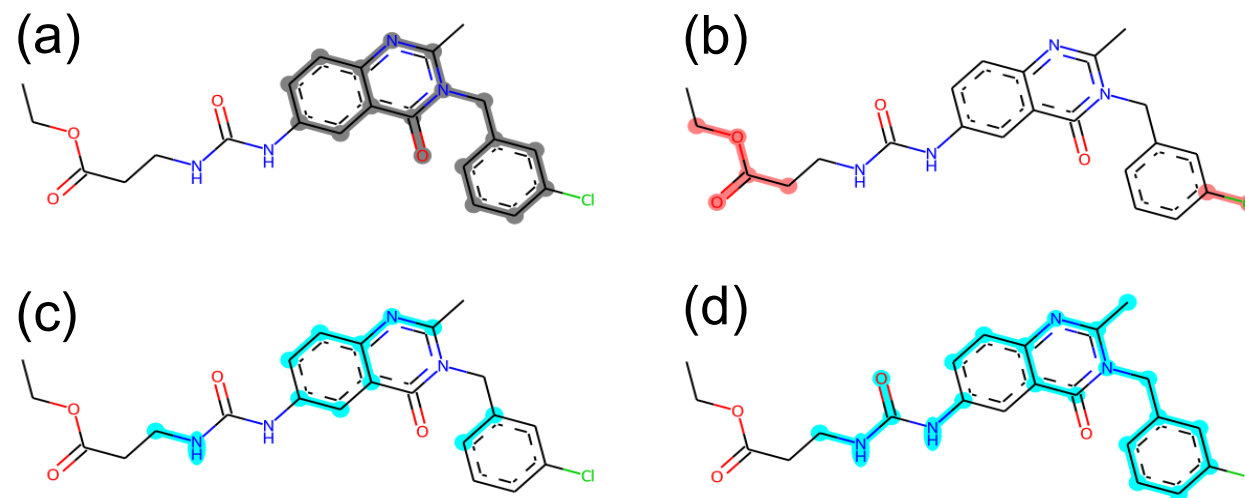
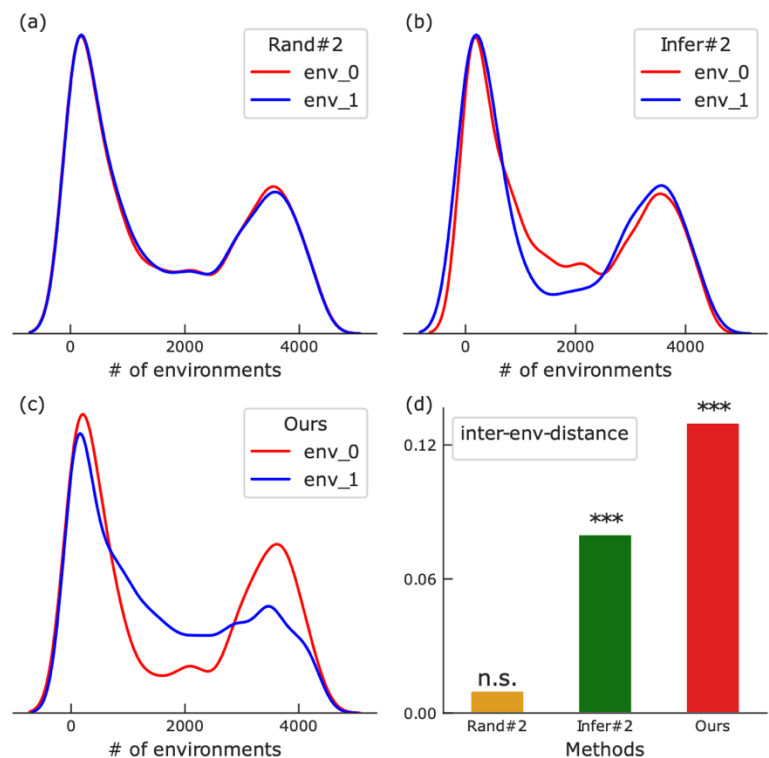
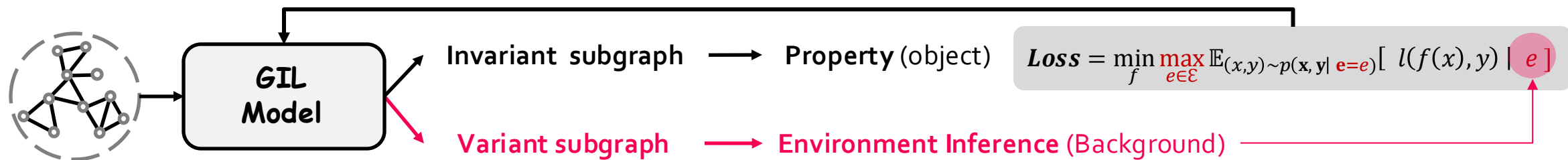
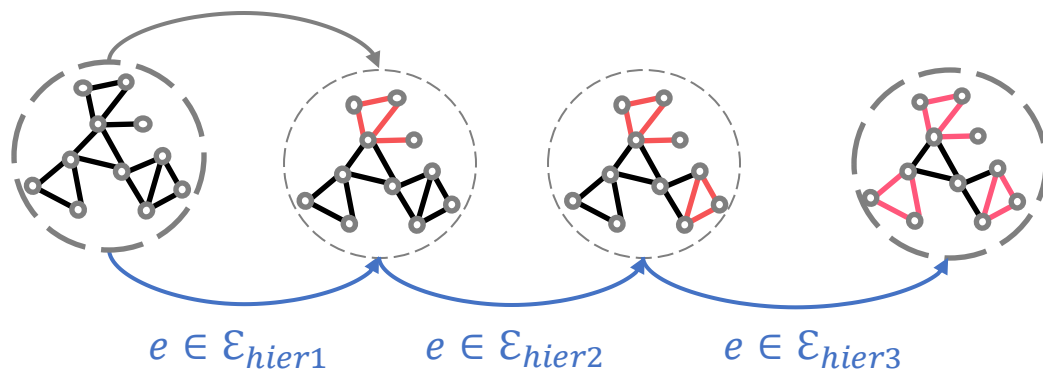


Figure 6. (a) Grey: Scaffold. (b) Red: Functional group (Chlorobenzene), and structural alert COC(=O)C from PubChem. (c, d) Blue: Learned variant subgraphs from the first and second hierarchies (the rest of the graph is considered learned invariant).

Conclusion



Existing methods EI in \mathcal{E}_{flat}



Our methods



Sufficiency

Efficiency

Interpretability

- Explicit Subgraph Modeling.
- Diverse environment inference.
- Env diversification loss instead of clustering methods.



Thank you for listening!

Results

- Effect of Hierarchical Semantic Environments (Table 4)
 - Environment Inference
 - Hierarchical Environment Inference
- Sensitive Analysis on Hierarchy (Table 5)

Table 4. Ablation studies on IC50-SCA and IC50-SIZE datasets.

| METHODS | IC50-SCA | IC50-SIZE |
|--------------------------------------|-------------------|-------------------|
| w/ env _{#non-infer(rand)=2} | 68.54±0.64 | 67.63±0.33 |
| w/ env _{#non-infer=real} | 68.77±0.72 | 67.60±0.32 |
| w/ env _{#infer=2} | 69.14±0.80 | 67.55±0.34 |
| w/ env _{#infer=real} | 69.08±0.64 | 67.74±0.13 |
| w/ env _{#hier-infer} (OURS) | 70.72±0.30 | 68.64±0.23 |

Table 5. Sensitivity analysis on generated environments.
(#e_p denotes the number of provided environments in datasets.)

| CONFIGURATIONS | IC50-SCA | IC50-SIZE |
|---|-------------------|-------------------|
| #real = [#e _p] | 69.08±0.64 | 67.60±0.32 |
| #env = [5] | 69.35±0.67 | 67.70±0.50 |
| #env = [2] | 69.14±0.80 | 67.73±0.64 |
| #env=[#e _p → #e _p /2 → 5] | 70.12±0.14 | 68.62±0.34 |
| #env=[#e _p → #e _p /2 → 2] | 70.72±0.30 | 68.64±0.23 |