

Evaluating Self-Supervised Learned Molecular Graph Representations

Anonymous Authors¹

Abstract

Because of data scarcity in real-world scenarios, obtaining pre-trained representations via self-supervised learning (SSL) has attracted increasing interest. Although various methods have been proposed, it is still under-explored what knowledge the networks learn from the pre-training tasks and how it relates to downstream properties. In this work, with an emphasis on chemical molecular graphs, we fill in this gap by devising a range of node-level, pair-level, and graph-level probe tasks to analyse the representations from pre-trained graph neural networks (GNNs). We empirically show that: 1. Pre-trained models have better downstream performance compared to randomly-initialised models due to their improved capability of capturing global topology and recognising substructures. 2. However, randomly initialised models outperform pre-trained models in terms of retaining local topology. Such information gradually disappears from the early layers to the last layers for pre-trained models.

1. Introduction

Self-Supervised Learning (SSL) pre-training has opened up the opportunity to effectively utilise vast amount of unlabelled data to improve downstream tasks where labels are limited. In natural language processing, language models like GPT-3 (Brown et al., 2020), Megatron (Shoeybi et al., 2019), and Gopher (Rae et al., 2021) can automatically re-discover the classical NLP pipeline in an interpretable and localisable way (Tenney et al., 2019). They can also achieve substantial improvements in a wide range of NLP tasks. In computer vision, self-supervised learning approaches such as contrastive learning (Chen et al., 2020c; He et al., 2020), bootstrapping (Grill et al., 2020) and masking (He et al., 2022) are shown to obtain competitive performance on widely-used benchmarks like ImageNet. DINO (Caron

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

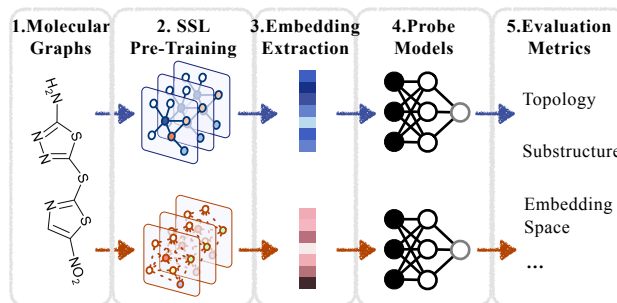


Figure 1. Overview of *GraphEval*. Given molecular graphs, we train GNNs to predict SSL proxy objectives. We then extract embeddings of (possibly unseen) graphs using pre-trained models, which form the inputs for probe models, trained and evaluated on the designed metrics.

et al., 2021a) shows that a self-supervised vision transformer (ViT) automatically learns class-specific features for unsupervised object segmentation.

Motivated by the successful applications of self-supervised learning, pre-training GNNs on unlabelled structured data has attracted increasing interest (Liu et al., 2021a; Xie et al., 2021). However, it is still under-explored what knowledge the networks learn during the pre-training and how it relates to downstream properties. In this work, with an emphasis on chemical molecules, we fill in this gap by devising: (1) a range of {node-, pair-, graph-} level metrics; (2) substructure detection; (3) embedding space characterisation, to analyse the representations from pre-trained GNNs. Our main insights are summarised as follows:

- Pre-trained representations are better at capturing global topological structure while losing the local information;
- Pre-trained models can well recognise molecular substructures that are correlated with properties;
- Pre-trained embedding space alleviates the issue of over-smoothing while the spectrum does not have obvious connections with embeddings’ quality.

2. Related Work

Graph SSL. Self-supervised learning methods for graphs are roughly categorised into contrastive and generative venues (Liu et al., 2021a;b; Wu et al., 2021; Xie et al., 2021). Contrastive graph SSL (Hu et al., 2020a; Sun et al., 2020;

You et al., 2020) applies contrastive learning to maximise the mutual information between augmented instances constructed from the same graph. Generative graph SSL (Hamilton et al., 2017; Hu et al., 2020a;b; Liu et al., 2018) forms the pretext task by reconstructing original graphs. A more recent trend in Graph SSL (Liu et al., 2022; Stärk et al., 2021) is to utilise domain knowledge, e.g., 3D information of molecular conformations, to help enhance the expressiveness of GNN. In this work, we focus on studying the transferable knowledge stored in the self-supervised learned molecular graph representations.

Probing Pre-trained Embeddings. Using probe models to study learned representations is a common practice to evaluate its quality. Probe models capture the intuition that good features should perform competitively in transfer tasks even with a shallow architecture. We review the related work applying probe models for natural language processing (Conneau & Kiela, 2018; Hendricks et al., 2021; Hewitt & Manning, 2019; Jawahar et al., 2019; Kassner & Schütze, 2020; Liu et al., 2019; Tenney et al., 2019; Wang et al., 2019), computer vision (Alain & Bengio, 2017; Caron et al., 2021b; Chen et al., 2020b; 2021; He et al., 2022; Li et al., 2021; Resnick et al., 2019; Wang et al., 2021), and biomedical science (Dohan et al., 2021; Elnaggar et al., 2021; Rao et al., 2019; Rives et al., 2021; Villegas-Morcillo et al., 2021). In natural language processing, pre-trained embeddings are shown to achieve competitive results on a wide range of tasks such as token labelling and parsing. In computer vision, self-supervised learned presentations can not only improve accuracy on downstream benchmarks such as ImageNet and CIFAR10, but also contain explicit semantic information (Caron et al., 2021b). In bioinformatics and biomedical science, self-supervised learning is able to learn biological structures and functions from massive unlabelled data. It has been shown that such learned embeddings are organised at a multi-scale level and can capture the information ranging from biochemical properties of amino acids to remote homological protein structures (Rives et al., 2021).

3. Preliminaries and Settings

We first introduce the basics of graphs and GNNs, then elaborate on the pre-training and probes.

Graph. A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ consists of a set of nodes \mathcal{V} and edges \mathcal{E} . In molecular graphs, nodes are atoms and edges are bonds. We use \mathbf{x}_u and \mathbf{x}_{uv} to denote the feature of node u and of the bond feature between nodes $[u, v]$, respectively. For notation simplicity, we use an adjacency matrix $\mathbf{A} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ to represent the graph, where $\mathbf{A}[u, v] \neq 0$ if the nodes (u, v) are connected.

GNN. There has been emerging research interest in exploring molecular graph representations (Corso et al., 2020;

Duvenaud et al., 2015; Gilmer et al., 2017; Liu et al., 2018; Yang et al., 2019). Graph neural networks are widely-adopted for encoding molecular graphs. A prototypical GNN uses messaging passing (Gilmer et al., 2017), where it updates atom-level representations based on their neighbourhoods. More specifically, let $\mathbf{h}_u^0 = \mathbf{x}_u$ be the input atom feature, we have:

$$\mathbf{m}_u^{t+1} = \sum_{v:\mathbf{A}[u,v] \neq 0} M_t(\mathbf{h}_u^t, \mathbf{h}_v^t, \mathbf{x}_{uv}), \quad \mathbf{h}_u^{t+1} = U_t(\mathbf{h}_u^t, \mathbf{m}_u^{t+1}) \quad (1)$$

where M_t and U_t are the message functions and vertex update functions, respectively. By repeating message passing for T steps, we can encode the information of the T -hop neighbourhood for each atom. We use a readout function R to pool node-level representations for graph-level prediction: $\hat{y} = R(\{\mathbf{h}_u^T | u \in \mathcal{V}\})$. In this work, we follow the research line of SSL on molecular graphs (Hu et al., 2020a; Liu et al., 2022; You et al., 2020) and adopt the Graph Isomorphism Network (GIN) (Xu et al., 2019) as the backbone model (modified in (Hu et al., 2020a) as to incorporate edge features during message passing).

Pre-Training. We use ten methods for Graph SSL, including EdgePred (Hamilton et al., 2017), InfoGraph (Sun et al., 2020), GPT-GNN (Hu et al., 2020b), AttrMask (Hu et al., 2020a), ContextPred (Hu et al., 2020a), G-Contextual, Motif (Rong et al., 2020), GraphCL (You et al., 2020), JOAO- $\{\cdot, v2\}$ (You et al., 2021) for pre-training. We follow the experimental settings and pre-training recipes reported in the original literature. For a fair comparison, we pre-train the same GIN model on the same data splits. Specifically, we randomly select 50k qualified molecules from the GEOM dataset (Axelrod & Gomez-Bombarelli, 2020). Once the pre-training finished, we extract the embeddings based on the saved weights and pass them to the probe tasks.

Probe. We use probe models (Liu et al., 2019) to study whether self-supervised learned representations encode helpful structural information about graphs. Concretely, we use a graph neural network to extract graph representations and train a shallow model to make predictions with these fixed node and graph embeddings. A common choice of the probe model (Hewitt & Liang, 2019) is either a linear projection or a multi-layer perceptron (MLP). We choose an MLP with one hidden layer to enable capturing the non-linear relations. We set the hidden size to 300 and apply the ReLU activation. We use scaffold splitting to split data into 80%/10%/10% for the training/validation/testing set. The training procedure runs for 100 epochs with a learning rate of $1e^{-3}$. We select the best model based on the validation set. All the results are averaged across three independent runs.

As follows, we show the effectiveness of SSL methods in downstream tasks and systematically study the knowledge

Table 1. Performance on molecular property predictions using probes, with (w/) or without (w/o) fine-tuning (FT). For each set of random/pre-trained embeddings, we report the ROC-AUC scores over 8 datasets consisting of 678 binary tasks, where the score of each task is averaged over three independent runs. We **bold** the best and underline the worst performance of each dataset.

FT	Random	AttrMask	GPT-GNN	InfoGraph	ContextPred	G-Contextual	G-Motif	GraphCL	JOAO	JOAOv2
w/o	<u>58.85</u>	62.18	61.43	61.94	59.58	64.63	62.59	63.00	60.99	62.31
w/	<u>67.21</u>	70.16	68.27	70.10	70.89	69.21	70.14	70.64	69.57	70.21

that the networks learn from the pre-training tasks:

- In Sec. 4, we evaluate SSL learned embeddings on molecular biochemical property, demonstrating that such substantial improvements with linear models and fine-tuning are not much relevant.
- In Sec. 5, we probe a wide range of structural and topological metrics based on the embeddings. We find that pre-trained embeddings are better at capturing global topological property, and randomised variants surprisingly outperform restoring local geometry.
- In Sec. 6, we demonstrate that pre-trained embeddings are better at predicting the counts of molecular substructures, *e.g.* allylic and benzene. We hypothesise that the superior performance of pre-trained embeddings for molecular biochemical property prediction comes from the fact that SSL pre-training help better capture the substructure existence (Alsentzer et al., 2020; Bouritsas et al., 2020).
- In Sec. 7, we characterise the embedding space of randomly initialised and SSL pre-trained methods. We analyse the embedding space’s spectrum, how the positive/negative embedding pairs align and distribute. We show that pre-training helps alleviate the over-smoothing issue, while the spectrum does not have obvious clues with embeddings’ performance.
- In Sec. 8, we dive deep on the randomised GNN models. By re-initialising parts of the pre-trained GNN models, analysing embeddings extracted at different stages of networks, we observe that different stages tend to capture different scales of graph structural information.

4. Biochemical Property Measure

We first use probe models to evaluate pre-trained embeddings on predicting molecular biochemical properties. Following previous graph SSL work (Hu et al., 2020a; You et al., 2020), we validate the quality of these embeddings on eight molecular datasets consisting of 678 binary property prediction tasks (Hu et al., 2021; Wu et al., 2018). As previously described in Sec. 3, for the setting of without fine-tuning (“w/o FT”), we update the probe models with fixed embeddings; with fine-tuning (“w FT”), both the pre-trained GNNs and the randomised probe models will be updated. We report the results in Table 1 and Table 9.

Results and Findings. As shown in Table 1, most of SSL pre-trained embeddings outperform the randomised peers both under fixed and non-fixed settings. Compared with fixed embeddings, tuning the pre-trained model weights will bring more substantial performance gains due to introducing more flexibility. However, in general, better performance at fixed embeddings does not accompany higher fine-tuning scores. For instance, embeddings pre-trained with “ContextPred” have the second-lowest score with fixed scenarios while perform the best after end-to-end fine-tuning. The correlation between the two sets of score rankings is 0.25, which questions the conventional approach’s rationale for evaluating the quality of learned embedding with linear models (He et al., 2022).

5. Topological Property Measure

We evaluate the pre-trained embeddings on metrics emphasising topological properties at multiple scales, which are based on the {node-, pair-, and graph-} level statistics. Many of these metrics are used as features in traditional machine learning pipelines on graphs prior to the advent of deep learning (Hamilton, 2020). We first provide descriptions of these metrics, then present results and findings.

Node-level statistics focus on local topological measures of a graph, where each node is accompanied with a metric value. They could be used as features in a node classification model (Hamilton, 2020).

- **Node Degree** (d_u) counts the number of edges incident to node u : $d_u = \sum_{v \in V} \mathbf{A}[u, v]$
- **Node Centrality** (e_u) represents a node’s importance, it is defined as a recurrence relation that is proportional to the average centrality of its neighbours:

$$e_u = \left(\sum_{v \in V} \mathbf{A}[u, v] e_v \right) / \lambda, \quad \forall u \in V \quad (2)$$

- **Clustering Coefficient** (c_u) measures how tightly clustered a node’s neighbourhood is:

$$c_u = (|(v_1, v_2) \in \mathcal{E} : v_1, v_2 \in \mathcal{N}(u)|) / d_u^2 \quad (3)$$

i.e. the proportion of closed triangles in neighbourhood (Watts & Strogatz, 1998).

Table 2. Performance on the topological metrics predictions. We report the mean square or the cross entropy loss (*i.e.*, the smaller the better), over all 8 downstream datasets. We **bold** the best and underline the worst performance of each metric. We have summarised the percentage where SSL pre-trained embeddings fail to outperform the random embeddings.

Metrics	Node			Pair			Graph			
	Pre-training	Degree	Centrality	Clustering	Link	Jaccard	Katz	Diameter	Connectivity	Cycle
–	0.001	1.199	0.297	31.05	1.879	<u>2.828</u>	<u>222.6</u>	<u>0.226</u>	<u>6.351</u>	<u>0.158</u>
AttrMask	0.015	1.307	0.424	32.23	2.029	2.634	164.7	0.178	6.075	0.102
GPT-GNN	3.032	1.380	0.505	41.44	<u>2.541</u>	2.374	178.8	0.247	9.222	0.166
InfoGraph	1.298	1.242	0.296	41.15	2.273	2.238	83.24	0.204	6.169	0.159
ContextPred	<u>5.498</u>	<u>1.626</u>	0.316	37.78	2.286	2.413	183.0	0.194	8.691	0.108
G-Motif	3.085	1.372	<u>0.531</u>	<u>51.83</u>	2.363	2.758	98.21	0.268	7.333	0.182
G-Contextual	0.036	1.242	0.403	33.55	1.773	2.660	113.6	0.170	5.330	0.045
GraphCL	0.854	1.110	0.461	34.97	1.863	2.271	89.79	0.226	6.191	0.152
JOAO	0.637	1.268	0.412	33.67	2.084	2.307	89.38	0.214	5.960	0.142
JOAOv2	0.591	1.272	0.463	32.81	2.054	2.340	88.27	0.217	5.964	0.148
SSL Worse	100%	89%	89%	100%	78%	0%	0%	0%	0%	0%

We use all the nodes from eight datasets, report the scores over eight test splits across multiple runs.

Graph-level statistics summarise global topology information and are helpful for tasks like graph classifications. We briefly describe their meanings and refer the formal definitions to (Hamilton, 2020).

- **Diameter:** maximum distance between the pair of nodes
- **Cycle Basis:** a set of simple cycles that forms a basis of the graph cycle space. It is a minimal set that allows every even-degree subgraph to be expressed as a symmetric difference of basis cycles.
- **Connectivity:** minimum number of elements (nodes or edges) that need to be removed to separate the remaining nodes into two or more isolated subgraphs.
- **Assortativity:** similarity of connections in the graph w.r.t the node degree, it is essentially the Pearson correlation coefficient of degree between pairs of linked nodes.

We use all the graphs from eight datasets, report the scores over eight test splits across multiple runs.

Pair-level statistics quantify the relationships between nodes. Since node and graph level statistics are not very useful for the tasks relied on relation modelling, we are interested in how well the pre-trained embeddings can capture the following pair-level metrics:

- **Link Prediction** tests whether two nodes are connected or not, given their embeddings and inner products. Based on the principle of *homophily*, it is expected that embeddings of connected nodes are more similar compared to disconnected pairs: $\mathbf{S}_{\text{Link}}[u, v, \mathbf{x}_u^T \mathbf{x}_v] = \mathbb{1}_{\mathcal{N}(u)}(v)$.

- **Jaccard Coefficient** seeks to quantify the overlap between neighbourhoods while minimising the biases induced by node degrees (Lü & Zhou, 2011): $\mathbf{S}_{\text{Jaccard}}[u, v] = |\mathcal{N}(u) \cap \mathcal{N}(v)| / |\mathcal{N}(u) \cup \mathcal{N}(v)|$
- **Katz Index** is a global overlap statistic, defined by the number of paths of all lengths between a pair of nodes: $\mathbf{S}_{\text{Katz}}[u, v] = \sum_{i=1}^{\infty} \beta^i \mathbf{A}^i[u, v]$, where $\beta \in \mathbb{R}^+$ is a pre-defined parameter controlling how much weight is given to short vs long paths. A small value ($\beta < 1$) down-weights the importance of long paths. Here we set $\beta = 1$, giving the paths of all lengths equal importance.

In experiments, we bootstrapped a fixed number of the node pairs (10k) from each dataset, report the test scores average over eight test splits across three runs.

Results and Findings. We report the results in Table 2. We observe that the randomised embeddings retain the local structural information well and outperform all the pre-trained embeddings. On the other hand, the pre-trained embeddings perform well when performing metrics related to the graph’s global topology. For pair-level statistics, randomised embeddings perform better when the metric itself is more about local structure, *e.g.* link prediction, and vice versa. We do not observe that there exists a dominant pre-training method that perform universally well w.r.t. other methods. There are some connections between the pre-training tasks and the performance on different metrics:

- Contextual proxy (*i.e.*, G-Contextual) is particularly helpful for Jaccard coefficient prediction because of the similarity of the pre-training objective and metric measure (neighbourhood overlap);

- Complicated design of augmentations (used in contrastive-based SSL, *i.e.* JOAO vs. GraphCL) do not bring substantial improvements in storing graph-level topological information.

6. Substructure Awareness Measure

Certain *substructures* usually reflect some properties at node and graph levels (Girvan & Newman, 2002). For instance, molecules containing benzene rings usually have similar physical (*e.g.* solvent) and chemical (*e.g.* aromaticity) properties (McMurry, 2014). On this basis, prediction (Alsentzer et al., 2020) and modelling (Bouritsas et al., 2020) of substructures have been proven effective for improving model expressiveness and downstream performance.

Molecular substructure. Instead of defining in an implicit or handcrafted manner, as in previous studies, a natural definition of substructure in molecules is the substituent or moiety that performs certain functions in chemical/biological reactions. Here we investigate 24 substructures which can be divided into three groups:

- **Rings:** Benzene, Beta lactams, Epoxide, Furan, Imidazole, Morpholine, Oxazole, Piperidine, Piperidine, Pyridine, Tetrazole, Thiazole, Thiophene
- **Functional Groups:** Amides, Amidine, Azo, Ether, Guanidine, Halogens, Hydroxylamine, Imide, Oxygens (including phenoxy), Urea
- **Redox Active Sites:** Allylic (excluding steroid dienone)

Each substructure might have unique effect on the downstream properties. For instance, forming with a simple cycle of atoms and bonds, a ring might lock particular atoms with distinct 3D structure therefore some of its stereochemistry properties such as chirality are determined, and chirality-aware modelling is proven beneficent in predicting molecular properties (Adams et al., 2022). We first apply “Cramér’s V” to measure how significant the substructures affect the molecular properties.

Cramér’s V quantifies the strength of the association between the molecular substructure counts (*i.e.*, chemical fragments) and their biochemical properties. It is defined as:

$$V = \sqrt{\chi^2 / (n \cdot \min(k - 1, r - 1))} = \sqrt{\chi^2 / n} \quad (r \equiv 2) \quad (4)$$

where n is the sample size, k and r are the total number of substructure counts and property categories (binary), respectively. The Chi-squared statistics χ^2 is then calculated as:

$$\chi^2 = \sum_{i,j} (n_{(i,j)} - n_{(i,\cdot)} \cdot n_{(\cdot,j)} / n)^2 / (n_{(i,\cdot)} \cdot n_{(\cdot,j)} / n) \quad (5)$$

Table 3. Cramér’s V between molecular substructure counts and binary properties, averaged over 678 binary property prediction tasks (*i.e.*, “Avg(Task)”) or eight downstream datasets (*i.e.*, “Avg(Data)”). We have also calculated the Pearson Rank Correlation (ρ) between embeddings’ performance on recognising the substructure and predicting properties.

Name	Type	Avg (Task)	Avg (Data)	ρ
allylic	Site	0.1144	0.1024	0.709
benzene	Ring	0.1630	0.1227	0.576
amide	Group	0.0881	0.1336	0.468
ether	Group	0.1034	0.1083	0.552
halogen	Group	0.1721	0.1086	0.515

Table 4. Performance on substructure detection. We **bold** the best and underline the worst performance of each substructure. It is clear to see that contrastive based method (GraphCL, JOAOv2) perform quite well in recognising these substructures. We provide detailed results in Appendix D.

Pre-training	allylic	amide	benzene	ether	halogen
–	3.516	<u>18.948</u>	<u>3.964</u>	6.071	<u>3.652</u>
AttrMask	3.371	12.932	2.860	4.958	1.192
GPT-GNN	2.808	15.736	2.938	5.932	2.912
InfoGraph	2.577	5.535	1.959	3.657	2.819
ContextPred	<u>4.386</u>	18.251	3.583	<u>7.045</u>	2.908
G-Motif	2.452	4.015	2.116	3.507	1.125
G-Contextual	2.196	5.938	1.926	2.900	0.759
GraphCL	2.088	3.922	1.722	3.766	0.798
JOAO	2.385	4.030	1.746	3.376	0.694
JOAOv2	2.122	3.865	1.773	3.388	0.695
SSL Worse	11%	0	0	11%	0

where $n_{(i,j)}$ is the total occurrence for the pair of (i, j) . Here i is the specific count of a certain substructure, and j represents the certain outcome of a molecular biochemical property. Cramér’s V value ranges from 0 to 1, representing the associated strength between two categorical variables.

Results and Findings. We calculate the Cramér’s V, and report the five substructures that are mostly correlated with downstream properties in Table 3. We report the detailed results in Table 11. We observe that certain molecular substructures are good indicators of their biochemical properties. Based on such facts, we train the probe models to predict the counts of substructures for all the molecules from the eight datasets. We report the test scores in Table 4. As noticed, all the pre-trained embeddings outperform random variants in terms of detecting the existence of substructures.

We also calculate the Pearson rank correlation ρ between the performance on downstream tasks and the performance on

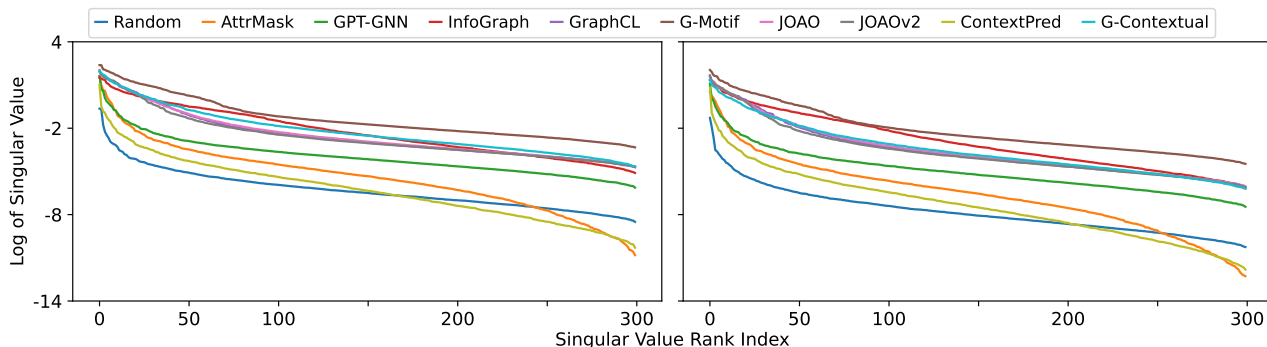


Figure 2. Spectrum of the SSL Pre-trained Embedding Space on BBBP. Left: Node; Right: Graph.

Table 5. Common classifiers trained based on the substructure counts for predicting molecular properties (ROC-AOC scores averaged over eight datasets). We utilised the conventional experimental setup in the sci-kit learn module. “Rand” and “SSL” represent the probe models trained on the randomised and GraphCL pre-trained embeddings, respectively.

Linear	RF	XGBoost	Probe (Rand)	Probe (SSL)
59.91	61.95	62.31	58.85	63.00

substructure detection of the SSL pre-trained embeddings. A strong positive correlation indicate that embeddings that are with better capability of detecting these substructures. Based on the observations of (1) molecular substructures are highly related with downstream biochemical properties; (2) embeddings that perform better in property predictions are usually with better substructure awareness; we conjecture that the performance gains from SSL pre-training might be from their capabilities of identifying graph substructures.

We find that: 1) substructure counts is highly correlated with the molecular properties; 2) the pre-trained embeddings are good at counting the substructures and predicting the properties. Consequently, we would like to measure that how well we can infer the properties solely based on the substructure counts.

How powerful are molecular substructure counters?

In question-answering systems, it has been found that the knowledge-aware graph modules may only carry out some simple reasoning such as counting (Wang et al., 2022). In GraphEval, we are interested in how the molecular substructure counters perform on the biochemical property predictions. We take the substructure counts as molecular descriptors to feed into classic methods, *e.g.*, linear classifier, random forest (RF), and XGBoost, which have been found (Jiang et al., 2021; Liu et al., 2018) to be effective in predicting molecular properties.

We report the averaged test ROC-AUC scores in Table 5.

Interestingly, these simple models trained on substructure counts achieve on par performance with SOTA 2D graph pre-trained embeddings. However, with more flexibility introduced by the end-to-end fine-tuning, the graph neural nets still maintain a margin of improvements ($\sim 7.7\%$). In retrospect to Table 1, we observe:

- with fixed pre-trained representation, GNN is comparative with substructure count descriptors + simple (linear) models;
- with fine-tuned representation, GNN perform much better than substructure counts.

Combining these two, we conjecture that GNN SSL pre-training strategies, especially contrastive-based, *e.g.* GraphCL and JOAO, are conducting something similar to substructure extraction/counting. However, it is not clear how fine-tuning pre-trained GNNs bring substantial improvements, we conjecture it might due to: (1) fine-tuning incorporate more information beyond substructure counting, such as pair/global topology; (2) GNN has larger model capacity which is born with more expressiveness.

7. Embedding Space Characterisation

Characterising the embedding space distribution helps to better understand notions of the amount of information preserved after being encoded. A prominent example is to measure the redundancy among embeddings (Wang & Isola, 2020; Zbontar et al., 2021). Here, we consider two related characterisations.

Spectrum (Dimensional Collapse). Following the similar protocols (Jing et al., 2022), we conduct the spectrum analysis of the learned embedding space. We first compute the covariance matrix $\mathbf{C} \in \mathbb{R}^{d \times d}$ based on all the embeddings \mathbf{z} :

$$\mathbf{C} = \frac{1}{N} \sum_{i=1}^N (\mathbf{z}_i - \bar{\mathbf{z}}) (\mathbf{z}_i - \bar{\mathbf{z}})^\top \quad (6)$$

After singular value decomposition:

$$\mathbf{C} = \mathbf{USV}^{\top} \quad (7)$$

where

$$\mathbf{S} = \text{diag}(\sigma^k) \quad (8)$$

we can plot the logarithm of the singular values ($\{\log(\sigma^k)\}$) in a descending order as in Fig. 2.

Dimension collapse, happening in visual contrastive SSL (Hua et al., 2021; Jing et al., 2022), refers to the problem of embedding vectors spanning a lower-dimensional subspace instead of the entire available embedding space. We observe that the graph pre-training methods that improve the downstream performance do not suffer from dimensional collapse. We conjecture the similar trends are due to the usage of the same baseline (GIN), therefore, sharing the same model capacity and expressiveness. We notice that a helpful pre-training will “lift” the spectrum, while a larger magnitude of the singular value (e.g. at the 50th/100th dimension) is usually with a better downstream performance.

Next we will analyse how the SSL pre-training affect the over-smoothness issue.

Alignment and Uniformity. Usually, similarities in the embedding spaces are expected to be conserved in the downstream properties. To examine this, we first clarify the meaning of positive and negative molecule pairs. Molecules that form positive pairs share *completely* identical biomedical properties provided by the dataset, while the negative pairs are the two that are entirely different. For instance, in Tox21 dataset, each molecule are labelled with 12 different biochemical properties (binary), only molecule pairs whose 12 properties are all the same are defined as positive.

We randomly select 10k positive and negative pairs from each dataset, calculate the cosine distance and plot the histogram in Fig. 3. We choose “AttrMask” and “GraphCL” to represent the generative and contrastive graph SSL pre-training methods.. The pre-trained embeddings formed two distinguishable distributions for positive/negative pairs, While the randomised embeddings do not. The positive/negative pair representations are indistinguishable, a phenomenon often referred to as *over-smoothing* (Chen et al., 2020a). This observation also explains why the randomised embeddings are less competitive in predicting downstream biochemical properties (Table 1) and capturing the graph statistics (Table 2). As SSL pre-training help alleviate the over-smoothing issues,, yet it is still worth investigating why the random embeddings perform quite well on capturing local structural information.

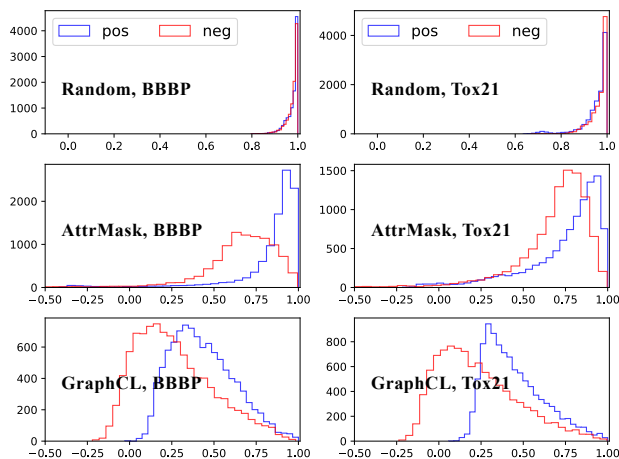


Figure 3. Cosine similarity between the positive and negative molecule pairs, sampled from BBBP and Tox21 datasets, respectively. The ‘pos’ and ‘neg’ stand for positive and negative molecule pairs, respectively

8. Randomised Features

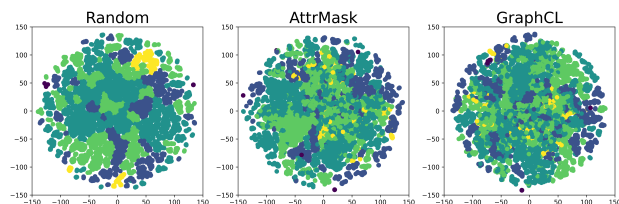


Figure 4. Visualisation on the graph embeddings using t-SNE, colored with the node degree, on the BBBP dataset.

Features extracted from randomly initialised networks are proved effective in applications such as face detection (Baek et al., 2021). While in GraphEval, we observe that randomised embeddings outperform all the pre-trained peers in recovering the node and (some) pair-level structural metrics. Such findings are validated by visualising the clusters using t-SNE in Fig. 4. In this section, we investigate how the pre-trained embeddings lose the local information. We re-initialised parts of the pre-trained weights to check how such perturbations result in the embeddings. We then compare the embeddings extracted from different stage layers of the network. We finally analyse how different message aggregation designs affect the extracted embeddings.

Different Module Re-initialisations. We start by observing how the extracted embeddings perform after re-initialising parts of the pre-trained GNN weights. The GIN architecture (Hu et al., 2020a; Xu et al., 2019) used in the Graph SSL research line consists of a node/edge feature embedding layer, and five subsequent GINConv layers and an output/readout layer. We use Glorot uniform initialisation

Evaluating Self-Supervised Learned Molecular Graph Representations

Table 6. Probe measures after re-initialising parts of weights of the pre-trained GNNs. We **bold** the best and underline the worst performance of each (pre-trained) embeddings on each structural metric. The digits in the “Re-Init” row represent the index/location of re-initialised GIN convolutional layers.

Metrics	AttrMask							GraphCL						
Re-Init	None	Embed	1	2	3	4	5	None	Embed	1	2	3	4	5
Node Degree	0.015	<u>0.263</u>	0.062	0.076	0.021	0.063	0.041	0.854	1.429	<u>1.439</u>	1.086	1.406	1.111	1.399
Graph Diameter	164.7	178.4	188.5	184.6	191.1	225.1	<u>229.7</u>	89.79	90.14	102.6	111.5	136.5	139.5	<u>163.9</u>

Table 7. Performance on the topological metrics predictions, based on the embeddings extracted from different layers. We report the mean square or the cross entropy loss (*i.e.*, the smaller the better), over test splits of 8 downstream datasets. We **bold** the best and underline the worst performance of each (pre-trained) embeddings on each metric.

Metrics	Random					AttrMask					GraphCL				
Stage (Layer No.)	1	2	3	4	5	1	2	3	4	5	1	2	3	4	5
Node Degree	<u>0.031</u>	0.013	0.007	0.001	0.001	0.009	<u>0.024</u>	0.019	0.013	0.015	0.016	0.361	0.426	0.583	<u>0.854</u>
Graph Diameter	184.0	<u>236.2</u>	167.4	184.0	222.5	<u>198.9</u>	181.9	184.6	171.7	164.7	93.34	<u>131.5</u>	97.43	95.19	89.79

for all network modules (Glorot & Bengio, 2010). We re-initialise the parts of pre-trained GNN weights with Glorot uniform initialisation (Glorot & Bengio, 2010), and extract the embeddings from the output layer, train probe models of these embeddings on the node degree (“ND”) and graph diameter (“GD”) metrics in Table 6. We observe:

- Re-initialising *any* layer might bring downside to the stored local topology info (node degree);
- Perturbing the late-stage (*i.e.* 4/5-th) layers would harm the learned hierarchical info (diameter).

Compared with random initialised networks, the later stage of the pre-trained message passing modules can better capture the global information. The embeddings extracted from the early stage contain more information on the local neighbourhood. We then study the outputs from different locations and different message passing schemes.

Different Stage Outputs. In table 7, we extract the embeddings from various stages of the random and pre-trained models and evaluate how they perform on predicting the structural metrics. We notice:

- Early-stage pre-trained embeddings outperform their randomised peers on local structural metrics, which might be because of better feature embedding modules;
- Such local information is gradually lost as the message passing modules stack, with more intention to represent the global graph signatures.

Different Aggregation Scheme. It has been reported (Dwivedi et al., 2020) that, isotropic aggregations

Table 8. Probe performance with different aggregation choices in the message passing. The ‘Sum’ aggregation is the default setting.

Metric	Aggr	Random	AttrMask	GraphCL
Degree	Max	31.69	12.87	12.36
	Sum	0.002	0.123	0.681
	Mean	0.003	0.064	1.381
Diameter	Max	156.2	181.1	87.93
	Sum	218.7	162.0	88.44
	Mean	208.2	287.2	117.2

are consistently better than the anisotropic counterparts on link predictions. Here we vary the aggregation methods in the graph convolution and report the performance in Table 8. We note that message passing designs have different effects with different pre-training paradigms. For instance, choosing the maximum messages seems to capture the most salient features of the entire graph (diameter) for random embeddings, but will also lose information about local structures (degree).

9. Discussion

In this work, we conduct a collection of probe tasks and analysis on evaluating the self-supervised learned graph embeddings. We conclude the performance gains introduced by the SSL pre-training come from a better awareness of global topology and substructures. The pre-trained message passing weights, help capture the hierarchical while hurdle the local information. A better design on the message passing module remains an open problem.

References

- Adams, K., Pattanaik, L., and Coley, C. W. Learning 3d representations of molecular chirality with invariance to bond rotations. In *ICLR*, 2022.
- Alain, G. and Bengio, Y. Understanding intermediate layers using linear classifier probes. In *ICLR Workshop*, 2017.
- Alsentzer, E., Finlayson, S., Li, M., and Zitnik, M. Subgraph neural networks. In *NeurIPS*, 2020.
- Axelrod, S. and Gomez-Bombarelli, R. Geom: Energy-annotated molecular conformations for property prediction and molecular generation. *arXiv:2006.05531*, 2020.
- Baek, S., Song, M., Jang, J., Kim, G., and Paik, S.-B. Face detection in untrained deep neural networks. *Nature Communications*, 2021.
- Bouritsas, G., Frasca, F., Zafeiriou, S., and Bronstein, M. M. Improving graph neural network expressivity via subgraph isomorphism counting. *arXiv:2006.09252*, 2020.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. In *NeurIPS*, 2020.
- Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., and Joulin, A. Emerging properties in self-supervised vision transformers. *arXiv:2104.14294*, 2021a.
- Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., and Joulin, A. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021b.
- Chen, D., Lin, Y., Li, W., Li, P., Zhou, J., and Sun, X. Measuring and relieving the over-smoothing problem for graph neural networks from the topological view. In *AAAI*, 2020a.
- Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Luan, D., and Sutskever, I. Generative pretraining from pixels. In *ICML*, 2020b.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *ICML*, 2020c.
- Chen, X., Xie, S., and He, K. An empirical study of training self-supervised vision transformers. *arXiv:2104.02057*, 2021.
- Conneau, A. and Kiela, D. Senteval: An evaluation toolkit for universal sentence representations. In *LREC*, 2018.
- Corso, G., Cavalleri, L., Beaini, D., Liò, P., and Veličković, P. Principal neighbourhood aggregation for graph nets. In *NeurIPS*, 2020.
- Dohan, D., Gane, A., Bileschi, M. L., Belanger, D., and Colwell, L. Improving protein function annotation via unsupervised pre-training: Robustness, efficiency, and insights. In *KDD*, 2021.
- Duvenaud, D., Maclaurin, D., Aguilera-Iparraguirre, J., Gómez-Bombarelli, R., Hirzel, T., Aspuru-Guzik, A., and Adams, R. P. Convolutional networks on graphs for learning molecular fingerprints. In *NeurIPS*, 2015.
- Dwivedi, V. P., Joshi, C. K., Laurent, T., Bengio, Y., and Bresson, X. Benchmarking graph neural networks. *arXiv:2003.00982*, 2020.
- Elnaggar, A., Heinzinger, M., Dallago, C., Rihawi, G., Wang, Y., Jones, L., Gibbs, T., Feher, T., Angerer, C., Steinegger, M., et al. Prottrans: towards cracking the language of life’s code through self-supervised deep learning and high performance computing. *IEEE PAMI*, 2021.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for quantum chemistry. In *ICML*, 2017.
- Girvan, M. and Newman, M. Community structure in social and biological networks. *PNAS*, 2002.
- Glorot, X. and Bengio, Y. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- Grill, J.-B., Strub, F., Altché, F., Tallec, C., Richemond, P. H., Buchatskaya, E., Doersch, C., Pires, B. A., Guo, Z. D., Azar, M. G., et al. Bootstrap your own latent: A new approach to self-supervised learning. In *NeurIPS*, 2020.
- Hamilton, W. L. Graph representation learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 14(3):1–159, 2020.
- Hamilton, W. L., Ying, R., and Leskovec, J. Inductive representation learning on large graphs. In *NeurIPS*, 2017.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. B. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. B. Masked autoencoders are scalable vision learners. In *CVPR*, 2022.
- Hendricks, L. A., Mellor, J., Schneider, R., Alayrac, J., and Nematzadeh, A. Decoupling the role of data, attention, and losses in multimodal transformers. *TACL*, 2021.
- Hewitt, J. and Liang, P. Designing and interpreting probes with control tasks. In *EMNLP*, 2019.
- Hewitt, J. and Manning, C. D. A structural probe for finding syntax in word representations. In *NAACL*, 2019.
- Hu, W., Liu, B., Gomes, J., Zitnik, M., Liang, P., Pande, V., and Leskovec, J. Strategies for pre-training graph neural networks. In *ICLR*, 2020a.
- Hu, W., Fey, M., Zitnik, M., Dong, Y., Ren, H., Liu, B., Catasta, M., and Leskovec, J. Open graph benchmark: Datasets for machine learning on graphs. In *NeurIPS*, 2021.
- Hu, Z., Dong, Y., Wang, K., Chang, K.-W., and Sun, Y. Gpt-gnn: Generative pre-training of graph neural networks. In *KDD*, 2020b.
- Hua, T., Wang, W., Xue, Z., Ren, S., Wang, Y., and Zhao, H. On feature decorrelation in self-supervised learning. In *ICCV*, 2021.

- Jawahar, G., Sagot, B., and Seddah, D. What does BERT learn about the structure of language? In *ACL*, 2019.
- Jiang, D., Wu, Z., Hsieh, C.-Y., Chen, G., Liao, B., Wang, Z., Shen, C., Cao, D., Wu, J., and Hou, T. Could graph neural networks learn better molecular representation for drug discovery? a comparison study of descriptor-based and graph-based models. *Journal of cheminformatics*, 2021.
- Jing, L., Vincent, P., LeCun, Y., and Tian, Y. Understanding dimensional collapse in contrastive self-supervised learning. In *ICLR*, 2022.
- Kassner, N. and Schütze, H. Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly. In *ACL*, 2020.
- Li, C., Yang, J., Zhang, P., Gao, M., Xiao, B., Dai, X., Yuan, L., and Gao, J. Efficient self-supervised vision transformers for representation learning. *arXiv preprint arXiv:2106.09785*, 2021.
- Liu, N. F., Gardner, M., Belinkov, Y., Peters, M. E., and Smith, N. A. Linguistic knowledge and transferability of contextual representations. In *NAACL*, 2019.
- Liu, S., Demirel, M. F., and Liang, Y. N-gram graph: Simple unsupervised representation for graphs, with applications to molecules. In *NeurIPS*, 2018.
- Liu, S., Wang, H., Liu, W., Lasenby, J., Guo, H., and Tang, J. Pre-training molecular graph representation with 3d geometry. In *ICLR*, 2022.
- Liu, X., Zhang, F., Hou, Z., Mian, L., Wang, Z., Zhang, J., and Tang, J. Self-supervised learning: Generative or contrastive. *IEEE TKDE*, 2021a.
- Liu, Y., Pan, S., Jin, M., Zhou, C., Xia, F., and Yu, P. S. Graph self-supervised learning: A survey. *arXiv:2103.00111*, 2021b.
- Lü, L. and Zhou, T. Link prediction in complex networks: A survey. *Physica A: statistical mechanics and its applications*, 2011.
- McMurry, J. E. *Organic chemistry with biological applications*. Cengage Learning, 2014.
- Rae, J. W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., Aslanides, J., Henderson, S., Ring, R., Young, S., et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv:2112.11446*, 2021.
- Rao, R., Bhattacharya, N., Thomas, N., Duan, Y., Chen, X., Canny, J., Abbeel, P., and Song, Y. S. Evaluating protein transfer learning with tape. In *NeurIPS*, 2019.
- Resnick, C., Zhan, Z., and Bruna, J. Probing the state of the art: A critical look at visual representation evaluation. *arXiv:1912.00215*, 2019.
- Rives, A., Meier, J., Sercu, T., Goyal, S., Lin, Z., Liu, J., Guo, D., Ott, M., Zitnick, C. L., Ma, J., et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *PNAS*, 2021.
- Rong, Y., Bian, Y., Xu, T., Xie, W., Wei, Y., Huang, W., and Huang, J. Self-supervised graph transformer on large-scale molecular data. In *NeurIPS*, 2020.
- Shoeybi, M., Patwary, M., Puri, R., LeGresley, P., Casper, J., and Catanzaro, B. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv:1909.08053*, 2019.
- Stärk, H., Beaini, D., Corso, G., Tossou, P., Dallago, C., Günnemann, S., and Liò, P. 3d infomax improves gnns for molecular property prediction. *arXiv:2110.04126*, 2021.
- Sun, F.-Y., Hoffmann, J., Verma, V., and Tang, J. Infograph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In *ICLR*, 2020.
- Tenney, I., Das, D., and Pavlick, E. BERT rediscovers the classical NLP pipeline. In *ACL*, 2019.
- Villegas-Morcillo, A., Makrodimitris, S., van Ham, R. C., Gomez, A. M., Sanchez, V., and Reinders, M. J. Unsupervised protein embeddings outperform hand-crafted sequence and structure features at predicting molecular function. *Bioinformatics*, 2021.
- Wang, A., Hula, J., Xia, P., Pappagari, R., McCoy, R. T., Patel, R., Kim, N., Tenney, I., Huang, Y., Yu, K., Jin, S., Chen, B., Durme, B. V., Grave, E., Pavlick, E., and Bowman, S. R. Can you tell me how to get past sesame street? sentence-level pretraining beyond language modeling. In *ACL*, 2019.
- Wang, H., Liu, Q., Yue, X., Lasenby, J., and Kusner, M. J. Unsupervised point cloud pre-training via occlusion completion. In *ICCV*, 2021.
- Wang, K., Zhang, Y., Yang, D., Song, L., and Qin, T. Gnn is a counter? revisiting gnn for question answering. In *ICLR*, 2022.
- Wang, T. and Isola, P. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *ICML*, 2020.
- Watts, D. J. and Strogatz, S. H. Collective dynamics of ‘small-world’ networks. *Nature*, 1998.
- Wu, L., Lin, H., Gao, Z., Tan, C., Li, S., et al. Self-supervised on graphs: Contrastive, generative, or predictive. *arXiv:2105.07342*, 2021.
- Wu, Z., Ramsundar, B., Feinberg, E. N., Gomes, J., Geniesse, C., Pappu, A. S., Leswing, K., and Pande, V. Moleculenet: a benchmark for molecular machine learning. *Chemical Science*, 2018.
- Xie, Y., Xu, Z., Zhang, J., Wang, Z., and Ji, S. Self-supervised learning of graph neural networks: A unified review. *arXiv:2102.10757*, 2021.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How powerful are graph neural networks? In *ICLR*, 2019.
- Yang, K., Swanson, K., Jin, W., Coley, C., Eiden, P., Gao, H., Guzman-Perez, A., Hopper, T., Kelley, B., Mathea, M., et al. Analyzing learned molecular representations for property prediction. *Journal of chemical information and modeling*, 2019.
- You, Y., Chen, T., Sui, Y., Chen, T., Wang, Z., and Shen, Y. Graph contrastive learning with augmentations. In *NeurIPS*, 2020.
- You, Y., Chen, T., Shen, Y., and Wang, Z. Graph contrastive learning automated. In *ICML*, 2021.
- Zbontar, J., Jing, L., Misra, I., LeCun, Y., and Deny, S. Barlow twins: Self-supervised learning via redundancy reduction. In *ICML*, 2021.

A. Probe Models

On the choice of architectures In principle, we want the probe models to be neither simple nor powerful. If the probe is too simple, then it won't be able to capture the information stored in the representations; on the other hand, a powerful probe itself will directly learn to precisely predict properties of our interests, since randomly initialised networks can already provide informative representations; therefore the predictive performance won't be a true reflection on the information contained in the representations. On such basis, we need to carefully choose the architecture of the probe models.

We first narrow down the search space to linear models (*i.e.*, mono-layer perceptron) and multi-layer perceptrons as they are widely used in previous studies (Alain & Bengio, 2017; Conneau & Kiela, 2018; Liu et al., 2019; Rao et al., 2019; Tenney et al., 2019). It then becomes simpler to quantify the capacity of the models based on their depth and width. To testify the capability of MLPs, we use embeddings pre-trained via AttrMask and GraphCL to predict two metrics, node-degree and graph diameter, on MLPs with various depth and width. We've tried that varying the number of hidden layers from 0 depths from 3, the dimensions from 100 to 1200. We find that: a linear model is too simple to capture the useful information while a three layers of MLPs will learn to predict by itself. Neither of them provide gave us much information and insights on the learned embeddings.

B. Pre-Training Setting

B.1. On the negative transfer of EdgePred

[to add descriptions and results](#)

B.2. Detailed results on the downstream property prediction

Table 9. Results for molecular property prediction tasks, with fixed embeddings. For each downstream task, we report the mean (and standard deviation) ROC-AUC of three seeds with scaffold splitting. The best and second best results are marked **bold** and **bold**, respectively. We have also reported the performance with end-to-end fine tuning. For fair comparison, we also train and report the average results of fine-tuning (Avg(FT)) downstream tasks.

Pre-training	BBBP \uparrow	Tox21 \uparrow	ToxCast \uparrow	Sider \uparrow	ClinTox \uparrow	MUV \uparrow	HIV \uparrow	Bace \uparrow	Avg \uparrow	Avg(FT) \uparrow
# Molecules	2,039	7,831	8,575	1,427	1,478	93,087	41,127	1,513	–	–
# Tasks	1	12	617	27	2	17	1	1	–	–
–	51.9(0.2)	65.6(0.4)	53.3(0.2)	57.6(0.1)	52.2(0.0)	53.7(0.4)	67.6(0.4)	69.1(0.2)	58.85	67.21
EdgePred	57.9(0.1)	53.2(0.6)	52.9(0.0)	48.6(0.4)	51.5(0.0)	53.9(0.2)	69.1(0.0)	67.6(0.3)	56.85	65.64
AttrMask	60.4(0.8)	66.7(0.7)	56.5(0.6)	60.0(0.4)	61.4(0.8)	46.2(0.7)	65.4(0.6)	80.8(0.4)	62.18	70.16
GPT-GNN	64.1(0.0)	62.1(0.1)	53.2(0.8)	58.4(0.1)	64.3(0.0)	46.5(0.1)	69.0(0.3)	73.8(0.3)	61.43	68.27
InfoGraph	66.3(0.6)	68.1(0.6)	58.4(0.6)	57.1(0.8)	66.3(0.6)	44.3(0.6)	70.2(0.6)	64.8(0.8)	61.94	70.10
ContextPred	54.6(0.2)	65.0(0.8)	56.5(0.1)	59.8(0.2)	52.1(0.2)	50.6(0.5)	65.3(0.1)	72.7(0.9)	59.58	70.89
G-Contextual	61.2(0.2)	69.7(2.3)	58.1(0.2)	61.7(1.9)	70.3(0.2)	44.1(1.6)	71.0(0.2)	80.9(1.9)	64.63	69.21
G-Motif	66.1(1.0)	65.6(0.7)	58.7(1.2)	58.7(1.1)	65.0(1.0)	39.2(0.7)	71.4(1.2)	75.9(1.1)	62.59	70.14
GraphCL	61.6(1.1)	67.5(0.8)	58.5(0.9)	57.6(0.8)	74.1(1.1)	44.0(0.8)	67.7(0.9)	73.1(0.8)	63.00	70.64
JOAO	63.8(0.7)	67.6(0.7)	57.1(0.7)	57.1(0.7)	59.2(0.7)	42.9(0.7)	69.4(0.7)	70.8(0.7)	60.99	69.57
JOAOv2	66.4(0.9)	68.2(0.8)	57.0(0.5)	59.1(0.7)	64.5(0.9)	47.4(0.8)	68.4(0.5)	67.4(0.7)	62.31	70.21
SSL Worse	10%	20%	10%	30%	20%	90%	20%	30%	10%	10%

B.3. In-Domain Generalisation

We pre-train on the same datasets that are used in the downstream tasks, it seems that in the local data regime, the in-domain pre-training does not provide useful performance gains.

Node Degree Scores on the test splits.

Evaluating Self-Supervised Learned Molecular Graph Representations

	BBBP	Tox21	ToxCast	Sider	ClinTox	MUV	HIV	Bace	Avg
Random 1	0.001	0.002	0.002	0.004	0.001	0.001	0.003	0.000	0.002
Random 2	0.002	0.001	0.003	0.003	0.001	0.001	0.003	0.000	0.002
Pre-trained in-domain datasets, AttrMask									
bbbp	0.240	0.444	0.455	0.162	0.140	1.301	0.987	0.199	0.491
tox21	1.795	2.220	2.815	1.139	1.247	8.288	6.513	1.179	3.149
toxcast	2.337	1.020	1.375	0.728	0.561	1.760	1.669	0.340	1.224
sider	0.162	0.353	0.325	0.124	0.122	0.650	0.663	0.085	0.311
clintox	0.195	0.291	0.298	0.239	0.131	0.180	0.492	0.058	0.235
muv	0.971	2.402	3.128	1.448	0.643	9.894	6.837	0.782	3.263
hiv	0.900	2.363	2.875	1.051	0.571	6.598	3.838	0.688	2.361
bace	0.107	0.268	0.408	0.182	0.116	0.468	0.617	0.100	0.283
Pre-trained in-domain datasets, GraphCL									
bbbp	0.267	0.221	0.290	0.542	0.363	0.097	0.088	0.128	0.250
tox21	0.420	0.622	0.456	0.835	0.141	0.190	0.193	1.166	0.503
toxcast	1.190	1.277	0.231	0.260	0.143	4.476	4.406	2.294	1.785
sider	0.138	0.111	0.049	0.058	0.312	0.368	0.225	0.071	0.166
clintox	0.071	0.059	1.669	1.204	1.316	0.084	0.030	6.563	1.374
muv	25.438	23.241	24.361	26.303	5.507	6.575	7.352	95.403	26.773
hiv	86.420	104.899	4.126	4.860	0.272	1.112	1.147	1.288	25.515
bace	0.247	0.212	0.257	0.266	6.807	4.660	4.478	0.199	2.141

B.4. On the Pre-Training Data

[add descriptions, update the tables and change the descriptions](#)

Table 10. Performance on the topological metrics predictions. We report the mean square or the cross entropy loss (*i.e.*, the smaller the better), over all 8 downstream datasets. We vary the amount of pre-training data (*i.e.*, GEOM) from 5k/10k/50k/100k/200k/330k(all) molecules. We use “↑” to represent the performance increase as the time increases, “↓” to represent the increase of the pre-train data hurdle the performance, “-” to represent there is no rules observed.

Metrics	Node			Pair			Graph			
	Degree	Centrality	Clustering	Link	Jaccord	Katz	Diameter	Connectivity	Cycle	Assortativity
AttrMask	↑	-	↑	-	↓	↑	-	↓	-	↓
GPT-GNN	↑	-	↑	-	-	↑	-	↑	-	↓
InfoGraph	↓	-	↑	-	↓	↑	-	-	-	-
ContextPred	↑	↑	-	↓	-	↑	-	↑	-	↓
G-Motif	↑	-	↑	-	↓	↑	-	-	-	↓
G-Contextual	↑	-	-	-	↓	↑	↑	-	-	-
GraphCL	-	-	↑	↓	-	↑	↑	-	↓	↑
JOAO	↑	-	↓	-	↓	↑	-	↑	-	↓
JOAOv2	↑	-	↓	-	↓	↑	-	↓	-	↓

C. On the embeddings from random initialised GNNs

We first analyse how the weights in the GNNs are initialised.

Edge Embedding layers ‘xavier uniform’, essentially are samples from uniform distribution

GNN layers essentially only have MLP weights (see [here](#)), same initialisation as Linear layers.

Linear layers (MLP) samples from uniform distribution for both weight and bias ([default in PyTorch](#))

Evaluating Self-Supervised Learned Molecular Graph Representations

Table 12. Substructure Detection, MSE Loss, Part I

Pre-training	allylic	amide	amidine	Azo	benzene	epoxide	ether	furan	guanido	halogen	imidazole	imide
–	3.516	18.948	0.370	0.216	3.964	0.146	6.071	0.536	0.217	3.652	0.531	0.417
EdgePred	5.854	25.126	0.417	0.251	9.804	0.154	12.838	0.705	0.249	5.214	0.597	0.578
AttrMask	3.371	12.932	0.364	0.174	2.860	0.123	4.958	0.242	0.213	1.192	0.459	0.319
GPT-GNN	2.808	15.736	0.360	0.156	2.938	0.141	5.932	0.407	0.220	2.912	0.483	0.474
InfoGraph	2.577	5.535	0.327	0.143	1.959	0.116	3.657	0.116	0.207	2.819	0.302	0.326
ContextPred	4.386	18.251	0.356	0.198	3.583	0.137	7.045	0.507	0.203	2.908	0.520	0.488
G-Motif	2.452	4.015	0.194	0.082	2.116	0.134	3.507	0.068	0.115	1.125	0.348	0.150
G-Contextual	2.196	5.938	0.172	0.067	1.926	0.114	2.900	0.084	0.121	0.759	0.249	0.160
GraphCL	2.088	3.922	0.262	0.124	1.722	0.127	3.766	0.074	0.167	0.798	0.274	0.192
JOAO	2.385	4.030	0.280	0.113	1.746	0.115	3.376	0.061	0.168	0.694	0.237	0.192
JOAOv2	2.122	3.865	0.262	0.112	1.773	0.123	3.388	0.068	0.125	0.695	0.244	0.197
SSL Worse	20%	10%	10%	10%	10%	10%	20%	10%	20%	10%	10%	20%

D. Substructure Detection

Descriptions For the descriptions of the chemical subgraphs, please refer to the [rdkit.Chem.Fragments](#).

Table 11. Cramér’s V between molecular substructure counts (categorical) and downstream properties (binary)

Pre-training	BBBP	Tox21	ToxCast	Sider	ClinTox	MUV	HIV	Bace	Avg (Task)	Avg (Data)
# Molecules	2,039	7,831	8,575	1,427	1,478	93,087	41,127	1,513	–	–
# Tasks	1	12	617	27	2	17	1	1	–	–
allylic	0.1602	0.1345	0.1156	0.1276	0.0935	0.0413	0.0280	0.1186	0.1144	0.1024
amide	0.2692	0.0490	0.0858	0.1841	0.1326	0.0235	0.0689	0.2556	0.0881	0.1336
amidine	0.0360	0.0291	0.0412	0.0323	0.0158	0.0117	0.0396	0.1328	0.0399	0.0423
azo	0.0400	0.0399	0.0393	0.0277	0.0123	0.0007	0.2082	-	0.0384	0.0526
benzene	0.1476	0.1632	0.1691	0.1149	0.1112	0.0289	0.1374	0.1091	0.1630	0.1227
epoxide	0.0273	0.0481	0.0449	0.0300	0.0049	0.0005	0.0086	-	0.0437	0.0235
ether	0.2314	0.0694	0.1060	0.1069	0.1023	0.0185	0.0498	0.1821	0.1034	0.1083
furan	0.0635	0.0257	0.0387	0.0227	0.0061	0.0311	0.0148	0.0135	0.0375	0.0270
guanido	0.0765	0.0201	0.0509	0.0715	0.0286	0.0057	0.0094	0.1088	0.0499	0.0464
halogen	0.1488	0.0849	0.1827	0.0773	0.0908	0.0143	0.0347	0.2353	0.1721	0.1086
imidazole	0.0601	0.0427	0.0492	0.0460	0.1212	0.0102	0.0398	0.1280	0.0483	0.0622
imide	0.0951	0.0246	0.0401	0.0428	0.0518	0.0094	0.0188	-	0.0392	0.0404
lactam	0.4263	0.0184	0.0116	0.0646	0.0543	0.0006	0.0048	-	0.0182	0.0830
morpholine	0.0512	0.0126	0.0343	0.0268	0.0425	0.0068	0.0101	0.0668	0.0329	0.0314
N_O	0.0438	0.0195	0.0467	0.0391	0.0709	0.0195	0.0144	0.0537	0.0452	0.0385
oxazole	0.0126	0.0184	0.0321	0.0359	0.0123	0.0079	0.0080	0.0364	0.0312	0.0205
piperidine	0.1450	0.0305	0.0844	0.0575	0.0418	0.0079	0.0226	0.0935	0.0803	0.0604
piperzine	0.0509	0.0214	0.0421	0.0776	0.0648	0.0111	0.0192	0.0063	0.0424	0.0367
pyridine	0.0598	0.0402	0.0549	0.0338	0.0833	0.0129	0.0300	0.1747	0.0529	0.0612
tetrazole	0.1161	0.0158	0.0251	0.0300	0.0286	0.0083	0.0123	0.0334	0.0247	0.0337
thiazole	0.1389	0.0521	0.0345	0.0445	0.0183	0.0118	0.0173	0.0539	0.0348	0.0464
thiophene	0.0356	0.0467	0.0472	0.0315	0.0113	0.0166	0.0081	0.0438	0.0456	0.0301
urea	0.0790	0.0236	0.0506	0.0471	0.0268	0.0079	0.0329	0.0516	0.0489	0.0399

Evaluating Self-Supervised Learned Molecular Graph Representations

Table 13. Substructure Detection, MSE Loss, Part II

Pre-training	lactam	morpholine	NO	oxazole	piperidine	piperzine	pyridine	tetrazole	thiazole	thiophene	urea
–	0.049	0.439	0.214	0.121	1.094	0.410	1.528	0.141	0.402	0.550	0.458
EdgePred	0.055	0.547	0.218	0.126	1.332	0.592	1.674	0.195	0.470	0.782	0.782
AttrMask	0.043	0.196	0.212	0.107	0.917	0.250	1.111	0.089	0.339	0.313	0.222
GPT-GNN	0.041	0.195	0.232	0.098	0.878	0.293	1.416	0.137	0.382	0.515	0.443
InfoGraph	0.041	0.121	0.241	0.054	0.813	0.206	0.655	0.075	0.170	0.180	0.174
ContextPred	0.043	0.298	0.223	0.120	1.097	0.387	1.334	0.146	0.379	0.538	0.538
G-Motif	0.048	0.153	0.212	0.031	0.768	0.219	0.670	0.034	0.081	0.092	0.092
G-Contextual	0.041	0.096	0.126	0.050	0.653	0.149	0.668	0.032	0.129	0.104	0.103
GraphCL	0.042	0.113	0.199	0.048	0.706	0.159	0.460	0.037	0.129	0.088	0.087
JOAO	0.043	0.096	0.194	0.040	0.738	0.149	0.434	0.032	0.121	0.077	0.076
JOAOv2	0.042	0.123	0.204	0.041	0.706	0.146	0.446	0.037	0.136	0.083	0.083
SSL Worse	10%	10%	40%	10%	20%	10%	10%	20%	10%	10%	20%

D.1. Cramer’s V

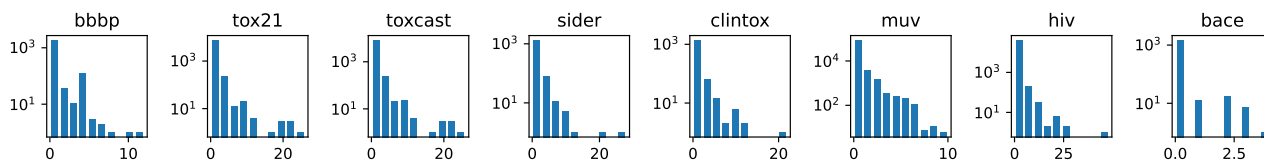
Table 14. Cramér’s V between molecular substructure counts (binary) and downstream properties (binary)

Pre-training	BBBP	Tox21	ToxCast	Sider	ClinTox	MUV	HIV	Bace	Avg (Task)	Avg (Data)
# Molecules	2,039	7,831	8,575	1,427	1,478	93,087	41,127	1,513	–	–
# Tasks	1	12	617	27	2	17	1	1	–	–
UNIFORM										
allylic	0.1602	0.1345	0.1156	0.1276	0.0935	0.0413	0.0280	0.1186	0.1144	0.1024
amide	0.2692	0.0490	0.0858	0.1841	0.1326	0.0235	0.0689	0.2556	0.0881	0.1336
amidine	0.0360	0.0291	0.0412	0.0323	0.0158	0.0117	0.0396	0.1328	0.0399	0.0423
azo	0.0400	0.0399	0.0393	0.0277	0.0123	0.0007	0.2082	-	0.0384	0.0526
benzene	0.1476	0.1632	0.1691	0.1149	0.1112	0.0289	0.1374	0.1091	0.1630	0.1227
epoxide	0.0273	0.0481	0.0449	0.0300	0.0049	0.0005	0.0086	-	0.0437	0.0235
ether	0.2314	0.0694	0.1060	0.1069	0.1023	0.0185	0.0498	0.1821	0.1034	0.1083
furan	0.0635	0.0257	0.0387	0.0227	0.0061	0.0311	0.0148	0.0135	0.0375	0.0270
guanido	0.0765	0.0201	0.0509	0.0715	0.0286	0.0057	0.0094	0.1088	0.0499	0.0464
halogen	0.1488	0.0849	0.1827	0.0773	0.0908	0.0143	0.0347	0.2353	0.1721	0.1086
imidazole	0.0601	0.0427	0.0492	0.0460	0.1212	0.0102	0.0398	0.1280	0.0483	0.0622
imide	0.0951	0.0246	0.0401	0.0428	0.0518	0.0094	0.0188	-	0.0392	0.0404
lactam	0.4263	0.0184	0.0116	0.0646	0.0543	0.0006	0.0048	-	0.0182	0.0830
morpholine	0.0512	0.0126	0.0343	0.0268	0.0425	0.0068	0.0101	0.0668	0.0329	0.0314
N_O	0.0438	0.0195	0.0467	0.0391	0.0709	0.0195	0.0144	0.0537	0.0452	0.0385
oxazole	0.0126	0.0184	0.0321	0.0359	0.0123	0.0079	0.0080	0.0364	0.0312	0.0205
piperidine	0.1450	0.0305	0.0844	0.0575	0.0418	0.0079	0.0226	0.0935	0.0803	0.0604
piperzine	0.0509	0.0214	0.0421	0.0776	0.0648	0.0111	0.0192	0.0063	0.0424	0.0367
pyridine	0.0598	0.0402	0.0549	0.0338	0.0833	0.0129	0.0300	0.1747	0.0529	0.0612
tetrazole	0.1161	0.0158	0.0251	0.0300	0.0286	0.0083	0.0123	0.0334	0.0247	0.0337
thiazole	0.1389	0.0521	0.0345	0.0445	0.0183	0.0118	0.0173	0.0539	0.0348	0.0464
thiophene	0.0356	0.0467	0.0472	0.0315	0.0113	0.0166	0.0081	0.0438	0.0456	0.0301
urea	0.0790	0.0236	0.0506	0.0471	0.0268	0.0079	0.0329	0.0516	0.0489	0.0399

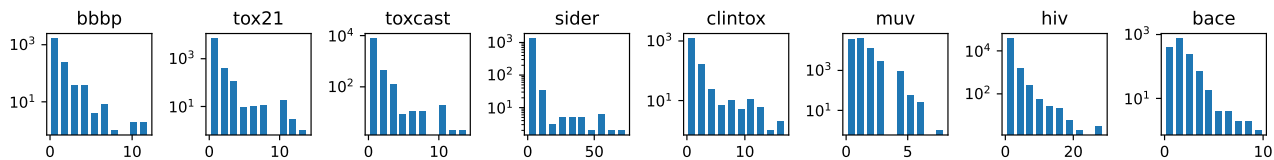
D.2. Pearson's chi-squared test, Substructure and Downstream Property

D.3. Distribution of Substructures

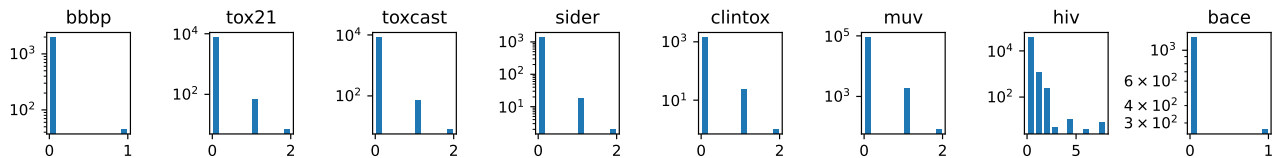
Allylic Oxide .



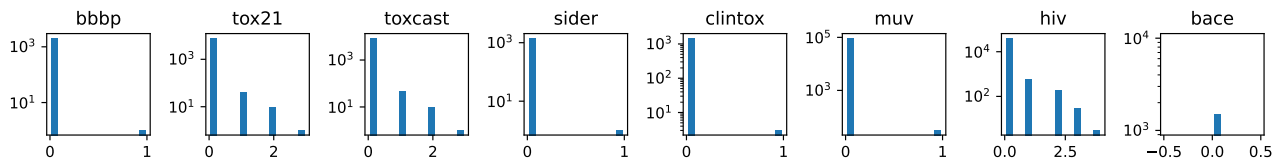
Amide .



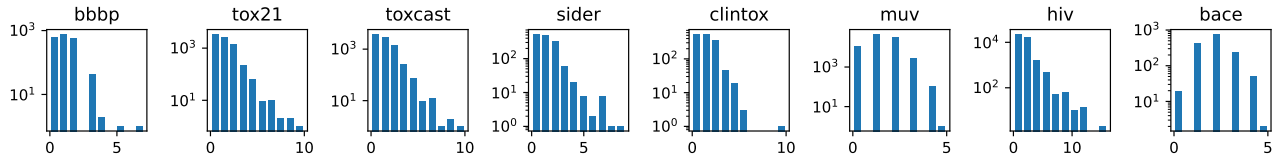
Amidine .



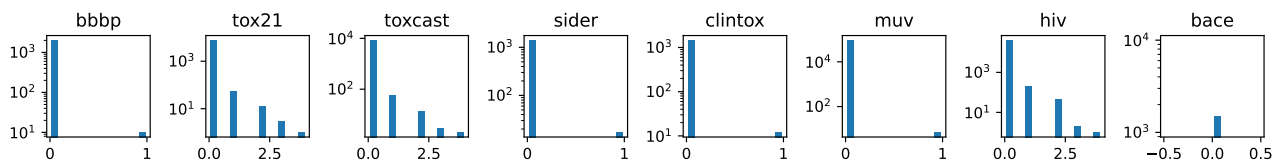
AZO .



Benzene .

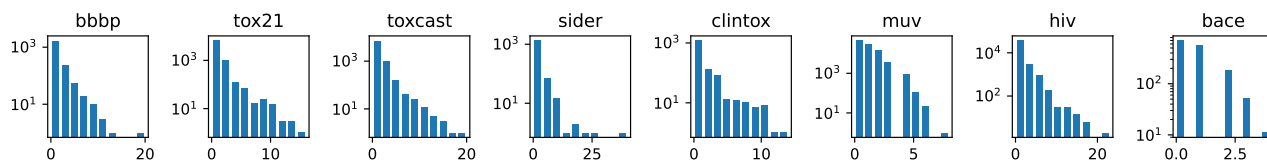


Epoxide .

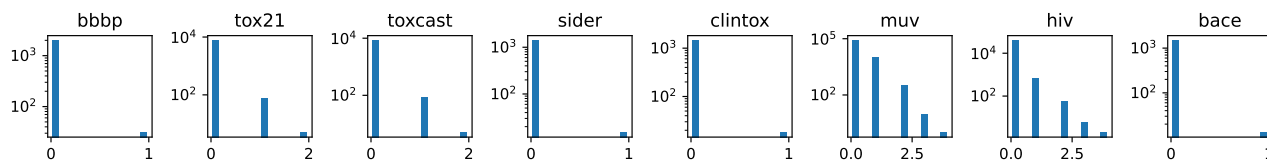


Evaluating Self-Supervised Learned Molecular Graph Representations

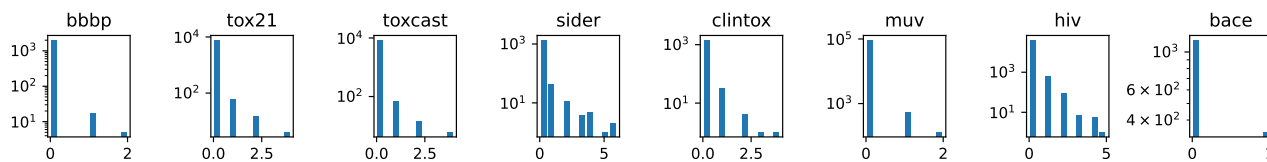
Ether



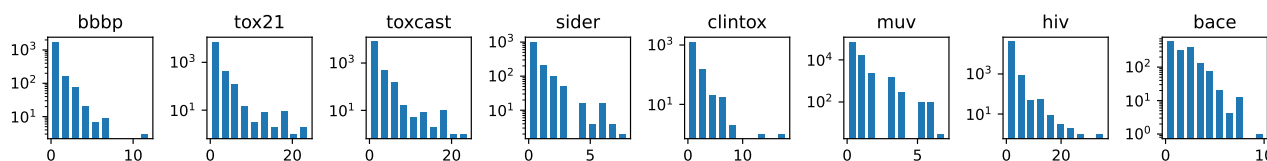
Furan



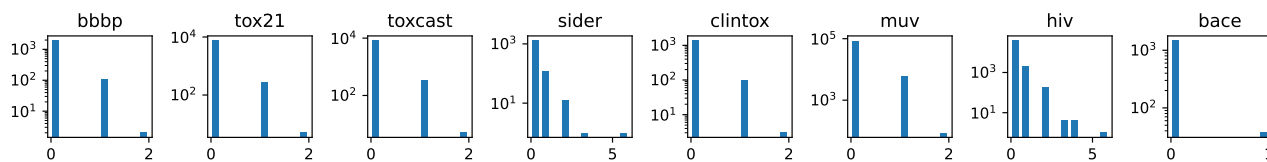
Guanido



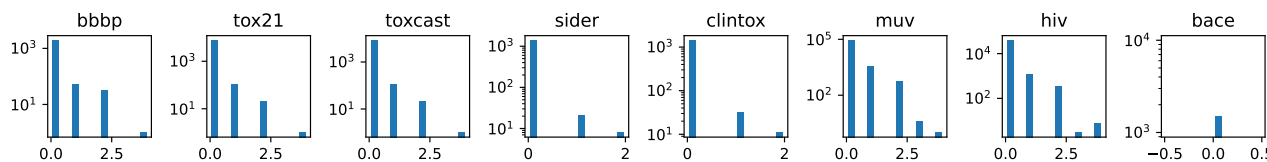
Halogen



Imidazole

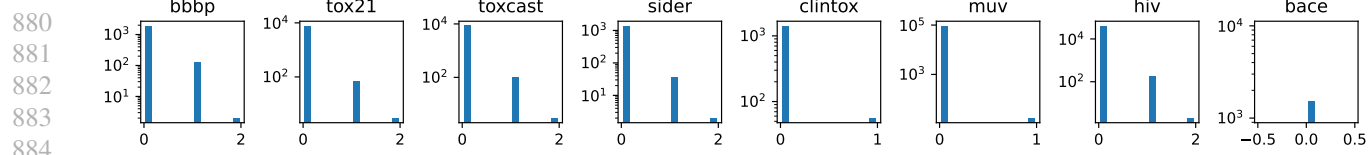


Imide

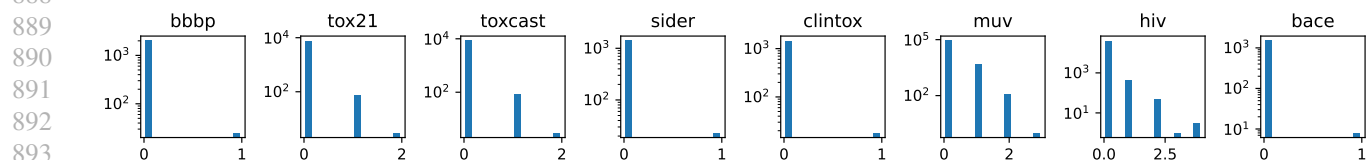


Lactam

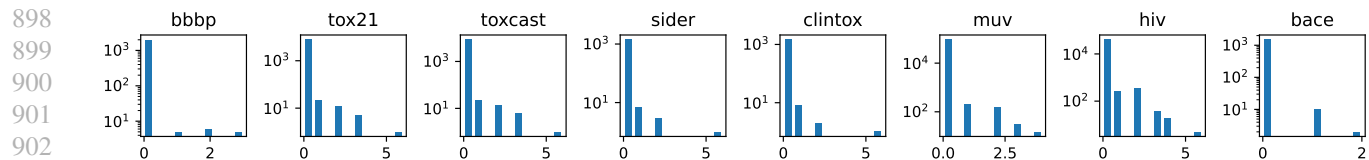
Evaluating Self-Supervised Learned Molecular Graph Representations



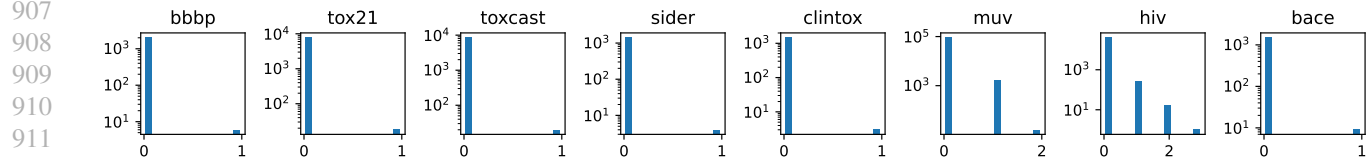
Morpholine



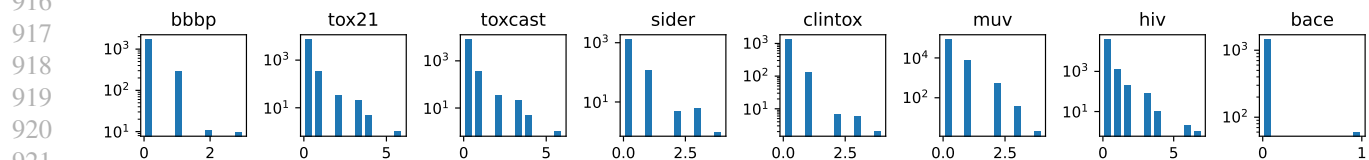
N_O



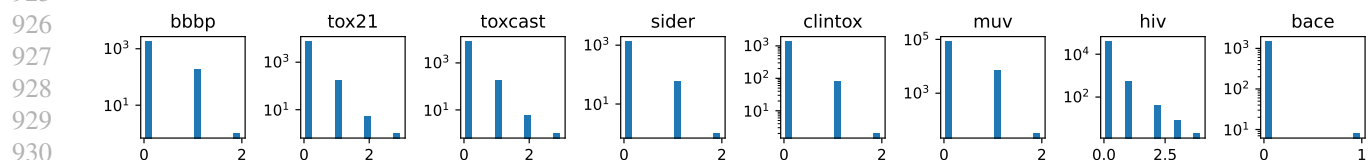
Oxazole



Piperidine



Piperzine

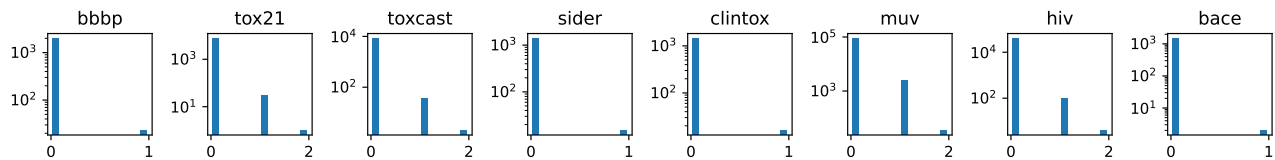
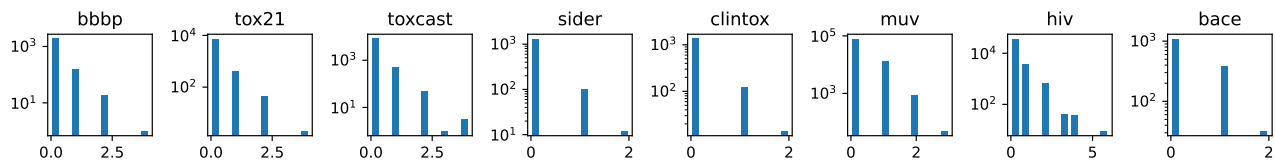


Pyridine

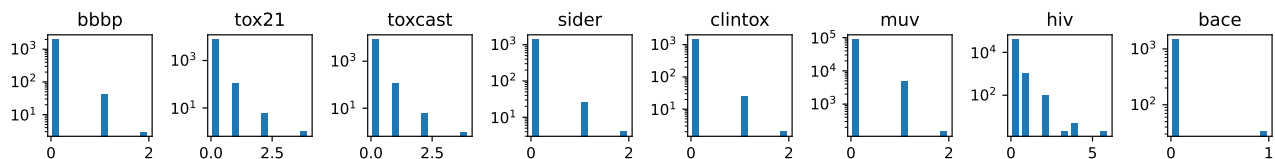
933
934

Evaluating Self-Supervised Learned Molecular Graph Representations

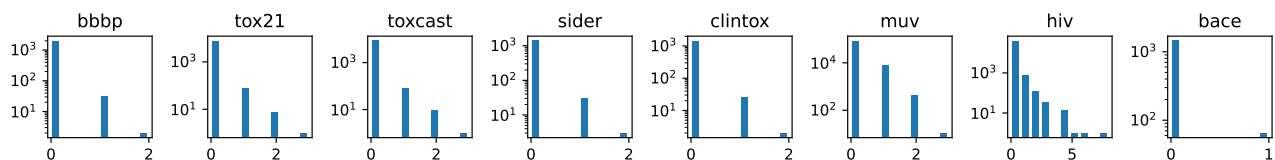
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989



Thiazole



Thiophene



Urea

E. Spectrum

We provide more examples on the spectrum analyses of the embedding space from different datasets.

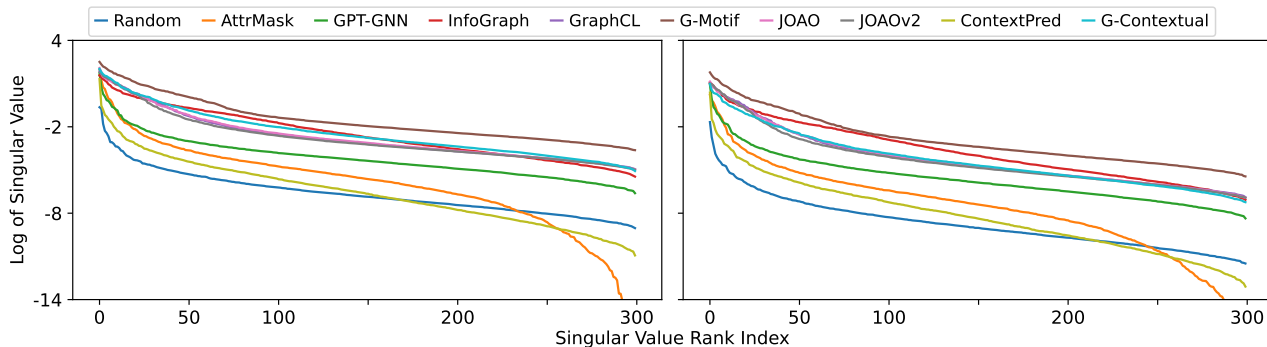


Figure 5. Spectrum of the SSL Pre-trained Embedding Space on Bace Dataset, Left: Node; Right: Graph.

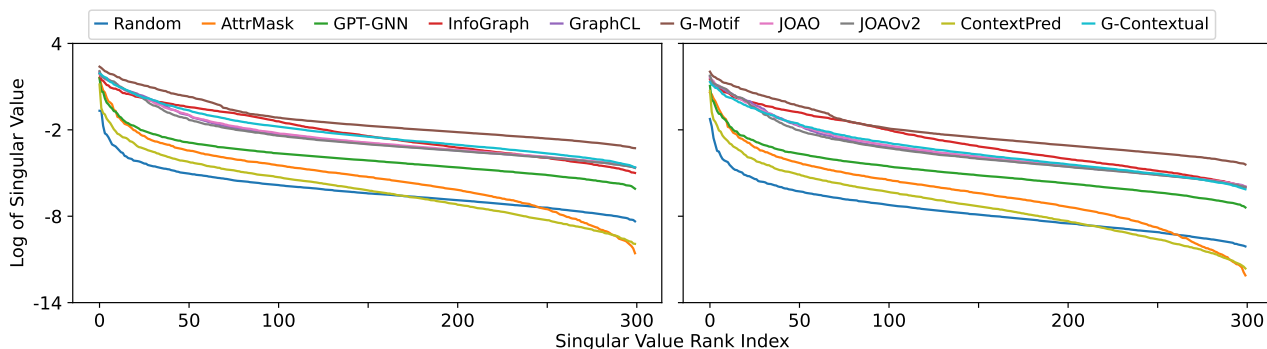


Figure 6. Spectrum of the SSL Pre-trained Embedding Space on Clintox Dataset, Left: Node; Right: Graph.

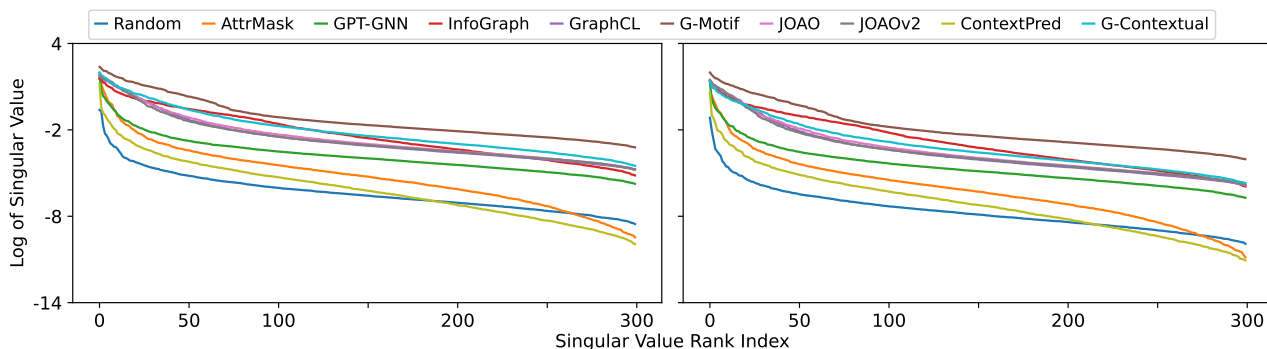


Figure 7. Spectrum of the SSL Pre-trained Embedding Space on HIV Dataset, Left: Node; Right: Graph.

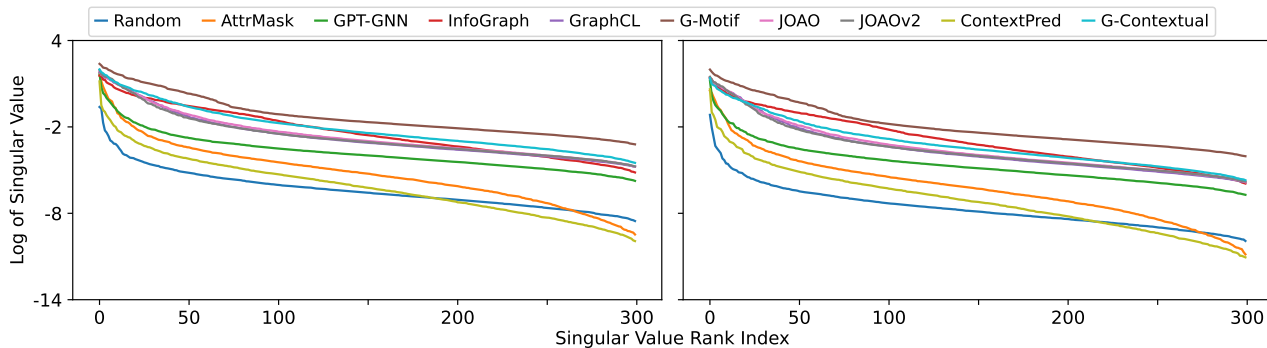


Figure 8. Spectrum of the SSL Pre-trained Embedding Space on HIV Dataset, Left: Node; Right: Graph.

Evaluating Self-Supervised Learned Molecular Graph Representations

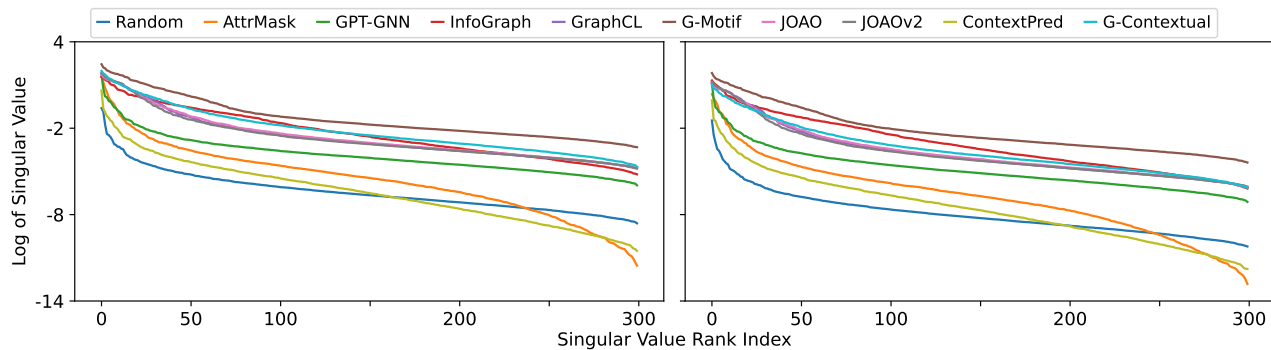


Figure 9. Spectrum of the SSL Pre-trained Embedding Space on MUV Dataset, Left: Node; Right: Graph.

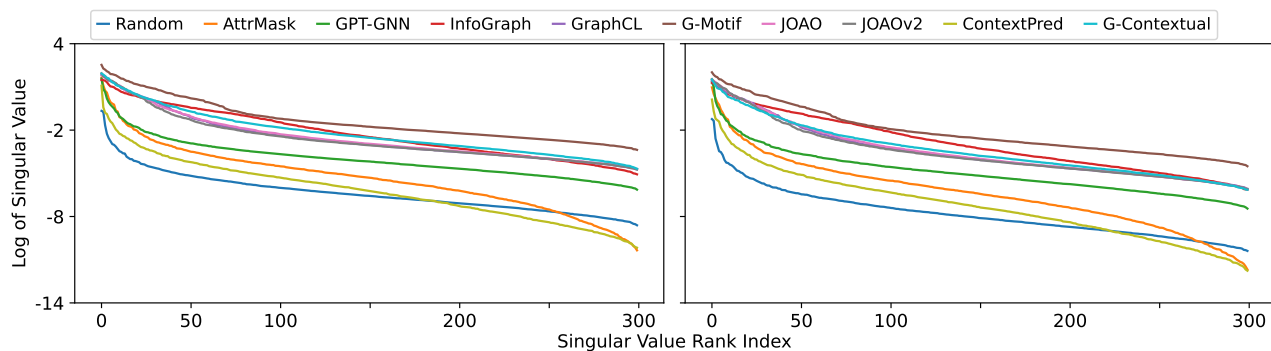


Figure 10. Spectrum of the SSL Pre-trained Embedding Space on Sider Dataset, Left: Node; Right: Graph.

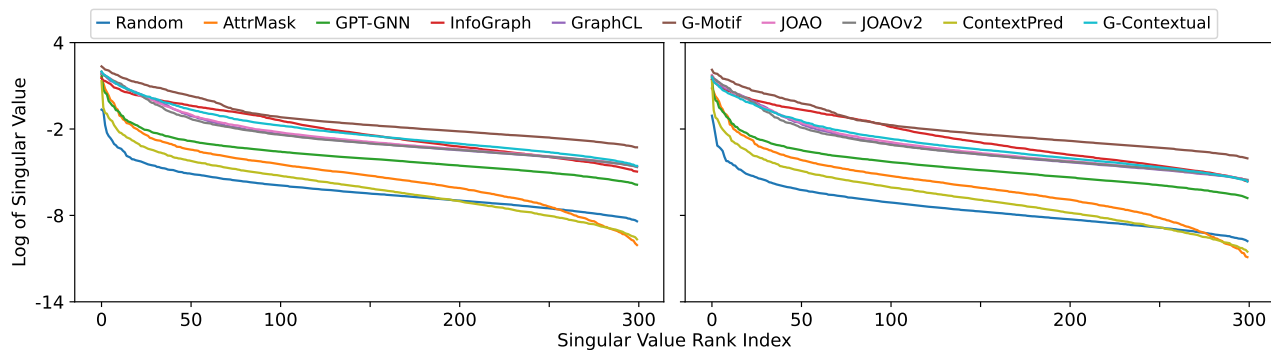


Figure 11. Spectrum of the SSL Pre-trained Embedding Space on Tox21 Dataset, Left: Node; Right: Graph.

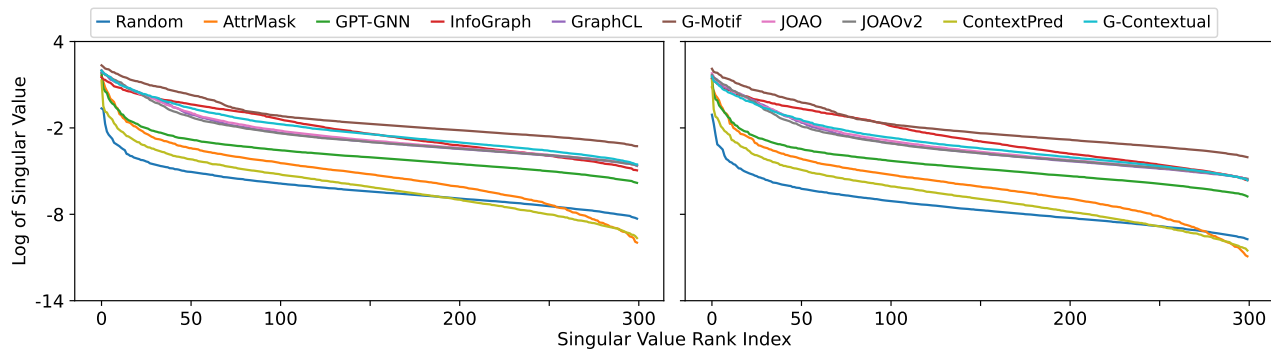


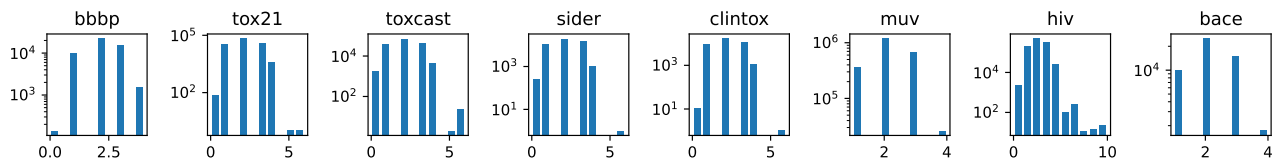
Figure 12. Spectrum of the SSL Pre-trained Embedding Space on Toxcast Dataset, Left: Node; Right: Graph.

F. Topological Metrics

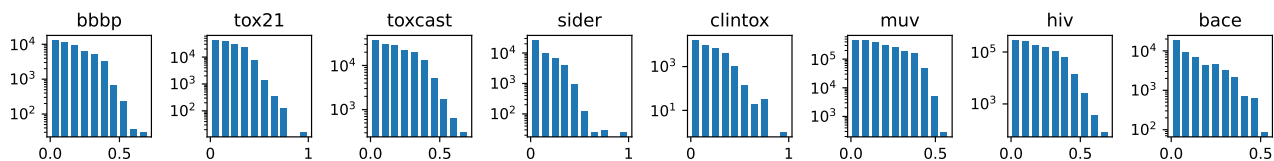
We first describe the graph structural metrics and motif substructures used in Sec. 6, we then visualise the histograms of the structural metrics and the substructure w.r.t. the downstream datasets. Note that the vertical axes sometimes are in the logarithm scale.

F.1. Distribution of Structural Metrics

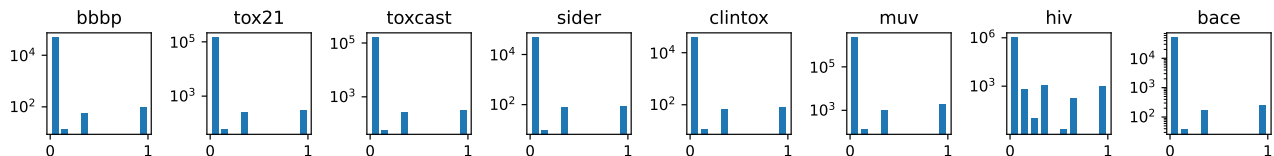
Node Degree



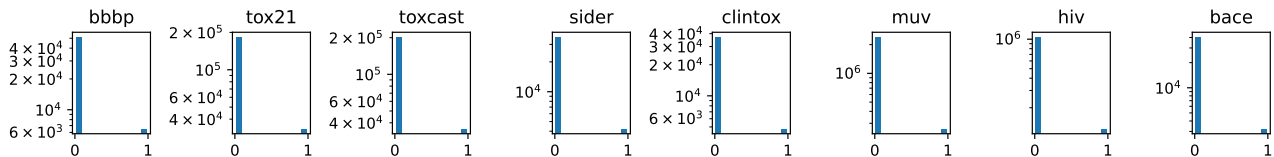
Node Centrality



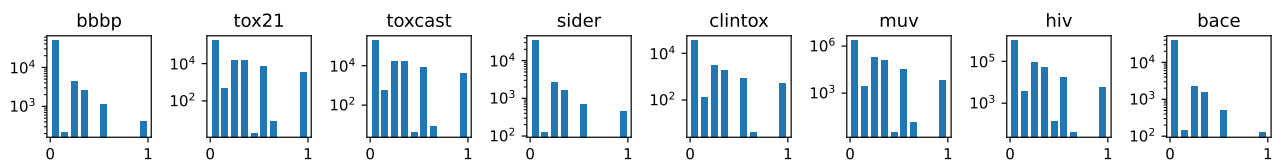
Node Clustering Coefficient



Link Prediction

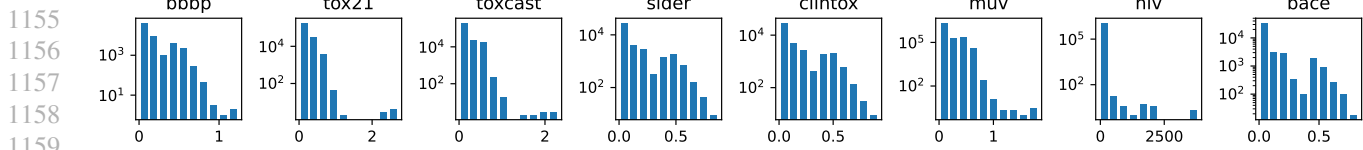


Jaccard Coefficient

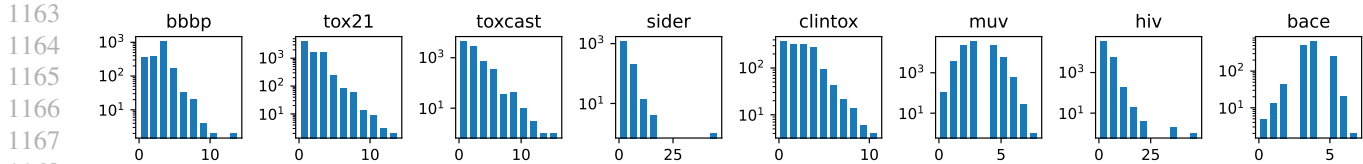


Katz Index

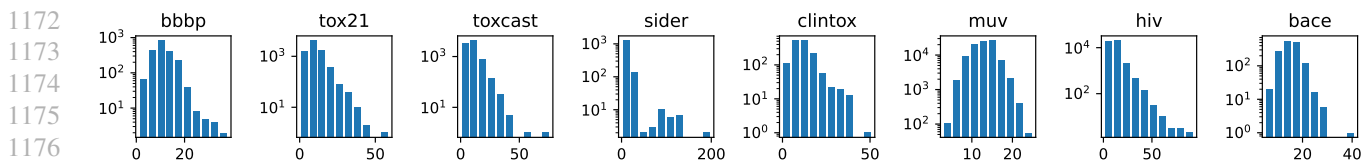
Evaluating Self-Supervised Learned Molecular Graph Representations



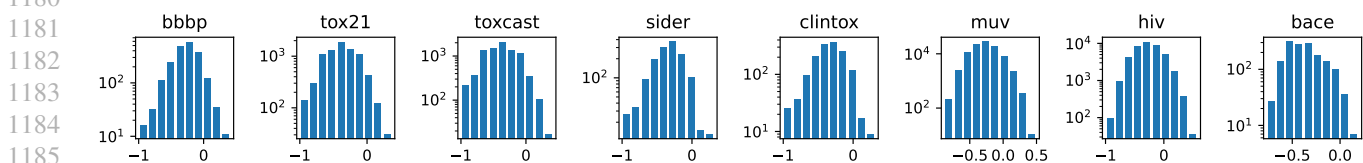
Cycle Basis



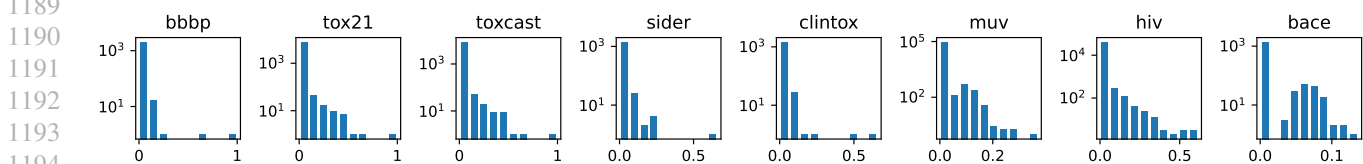
Graph Diameter



Assortativity Coefficient



Average Clustering Coefficient



F.2. Descriptions on the Graph-Level Structural Metrics

1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209

We refer to “Graph Representation Learning” written by WL Hamilton the for the description of graph-level structural metrics used in this study.