

# PREDICTING NETWORK MOTIF FINGERPRINTS WITH GRAPH NEURAL NETWORKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Graph Neural Networks (GNNs) are a predominant method for graph representation learning. However, beyond subgraph frequency estimation, their application to network motif prediction remains underexplored, with no established benchmarks in the literature. We propose to address this problem, framing motif estimation as an extension of subgraph frequency estimation. Our approach formulates motif estimation as a multitarget regression problem, optimising for interpretability and improving stability and scalability on large graphs. We validate our method using a large synthetic dataset generated by graph generators that mimic real-world data, and further test it on real-world graphs. Our experiments reveal that 1-WL limited models trained on synthetic data struggle to predict accurately motif profiles of real-world networks. However, apart from their reasonable performance within synthetic data, they can generalise to approximate the graph generation processes of real-world networks by comparing their predicted motif profiles with the ones originating from synthetic data. This first study on GNN-based motif estimation sets a benchmark and should open pathways for further developing the connection between motif profiles and subgraph frequency from a graph representation learning perspective.

## 1 INTRODUCTION

A structure is called a network motif when its recurring occurrence is not solely explained by randomness. These structures are extremely powerful tools for understanding complex networks. Understanding what substructures are relevant and not relevant to a graph can help understand the fundamental organisational principles behind it. This understanding enhances theoretical knowledge of network structure and function but also has practical implications in various fields, particularly in Biology. For instance, the feed-forward loop has been identified as a crucial functional pattern in many real biological networks of gene regulation (Mangan & Alon, 2003). It has also been discovered that motifs enable efficient communication and fault-tolerance across transcriptional networks (Roy et al., 2020). Furthermore, the related concept of graphlet degree distribution – a generalisation of degree distribution to higher-order structures – has been used to understand what is a good network model for protein-protein interactions (Pržulj, 2007).

Discovering a motif entails counting the number of occurrences of the desired structure, both in the network in study and in a set of control networks to understand its significance. However, this process is a very hard computational task. Just determining if a subgraph exists in a larger network (subgraph isomorphism) is a NP-complete problem (Cook, 1971). Even though methods to perform an analysis based on motifs exist (Ribeiro et al., 2021), they have high temporal complexity, rendering them intractable for very large networks. Furthermore, methods that rely on machine learning to address this problem typically do not give very interpretable results, a critical concept when doing analysis based on the relative importance of substructures.

**Present Work.** We aim to design a method for motif finding, leveraging a novel formulation that hinges primarily on reworking the target task to something else other than direct substructure counting. Our approach focuses on providing highly interpretable scores, ensuring the possibility of further insight into the conclusions obtained. Additionally, our method is robust and versatile, capable of operating effectively on graphs of any size. Knowing how difficult it is to obtain a high volume of real-world graph datasets that have both high quality and variety, we create a large synthetic dataset,

employing a myriad of generators and using it as the training data to assess the efficiency of our formulation. This setup leads to the conceptualisation of the following research question. Can Message Passing Neural Networks (MPNNs) under a different problem formulation than graph counting, be enough to accurately predict motifs of real-world graphs when trained on synthetic data?

**Key Contributions:** Our key contributions can be summarised as follows: **(1)** We show that the difficulty of motif discovery with MPNNs can be manipulated through different formulations of the target variable. Different formulations pertain to different concepts of motifs. Hence, depending on the the concept used, motif estimation does not have to follow the limitations from the literature regarding subgraph counting with MPNNs; **(2)** We make available a large diverse synthetic dataset in terms of graph topology and motif significance-profile generated with 23 synthetic generators. We also present a collection of more than 100 real-world networks and their motif significance-profile; **(3)** Our experiments show that MPNNs trained under the adopted motif concept with synthetic data can predict the significance-profile of synthetic data with a solid accuracy. Furthermore, we show they cannot generalise well for real-world data, showing a gap between these types of networks.

## 2 PRELIMINARIES

Let a graph  $\mathcal{G} = (\mathbb{V}_{\mathcal{G}}, \mathbb{L}_{\mathcal{G}}, \mathbf{X})$  where  $\mathbb{V}_{\mathcal{G}}$  denotes the vertex set of  $\mathcal{G}$ ,  $\mathbb{L}_{\mathcal{G}} \subseteq \mathbb{V}_{\mathcal{G}} \times \mathbb{V}_{\mathcal{G}}$  the edge set of  $\mathcal{G}$ , and  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$  the vertex features such that  $\forall v \in \mathbb{V}_{\mathcal{G}}, \mathbf{x} \in \mathbb{R}^{d_2}$ . Let all edges be undirected such  $(u, v) \in \mathbb{L}_{\mathcal{G}} \Leftrightarrow (v, u) \in \mathbb{L}_{\mathcal{G}}$ . Let  $\mathcal{H}$  be a subgraph of  $\mathcal{G}$  if and only if  $\mathbb{H}_{\mathcal{H}} \subseteq \mathbb{V}_{\mathcal{G}} \wedge \mathbb{L}_{\mathcal{H}} \subseteq \mathbb{L}_{\mathcal{G}}$ , such that exists an injective homomorphism given by the injective function  $f : \mathbb{V}_{\mathcal{H}} \mapsto \mathbb{V}_{\mathcal{G}}$  such that  $(v, u) \in \mathbb{L}_{\mathcal{H}} \Rightarrow (f(v), f(u)) \in \mathbb{L}_{\mathcal{G}}$ . If  $f$  is bijective and  $f^{-1}$  is an homomorphism (injective by construction) the relation is an isomorphism and the subgraph induced.

## 3 RELATED WORK

In order to discover motifs, we must define three steps: (1) What is the set of graphs,  $S_{\mathcal{G}}$ , that we admit as candidates for motifs; (2) What method is used to count the occurrences of graphs of  $S_{\mathcal{G}}$  in the graph of interest  $\mathcal{G}$ ; (3) How is the significance of the obtained counts calculated.

### 3.1 STEP ONE - HOW TO DEFINE THE SET OF GRAPHS USED

In this step,  $S_{\mathcal{G}}$  is typically defined *a priori*. This method is the most widely used (Milo et al., 2004a;b; Shen-Orr et al., 2002), and in most common cases, the selection of graphs used are ones known to be important to the area of the work in question (Shen-Orr et al., 2002; Alon, 2007).

Defining  $S_{\mathcal{G}}$  *a priori* is frequent for “non-machine-learning” techniques, but it is also common in machine-learning ones (Rong et al., 2020; Ying et al., 2020). However, when using techniques based on machine-learning, it is easier to create a task that can infer structures in  $\mathcal{G}$  than when using non-machine learning approaches. Hence, motif discovery can be modulated as the task of finding the best set of graphs that are motifs according to a defined graph metric. That is, discover what graphs exhibit a certain criteria in order to be considered motifs (Bénichou et al., 2023; Zhang et al., 2020).

### 3.2 STEP TWO - HOW TO COUNT SUBGRAPHS

**Non-GNN Methods.** Numerous methods exist for approximating subgraph counts, eschewing dependence on Graph Neural Networks (GNNs) or any machine learning techniques. We refer the interested reader to Ribeiro et al. (2021) for a survey of these methods.

**GNN Methods.** Counting occurrences of a graph  $\mathcal{G}$  in another graph  $\mathcal{H}$  using GNNs was first introduced by Chen et al. (2020). This work introduced significant limitations of what substructures MPNNs can count. Subsequent works have refined MPNN-like models to be more expressive, allowing them to have guarantees of being able to count occurrences of more graphs. One branch of such models is known as node-rooted Subgraph GNNs. These will extract a  $k$ -hop neighbourhood for each node in the graph to be studied. Since they act per node and add a feature to the node that induced each subgraph, they are called node-rooted Subgraph GNNs.

108 These architectures, with models as powerful as 1-WL as the backbone, are strictly more powerful  
 109 than maximum powerful MPNNs but are less powerful than the 3-WL test (Frasca et al., 2022; Yan  
 110 et al., 2023; Zhang et al., 2023). Hence, they have limitations regarding the type of structures that  
 111 they can count. Huang et al. (2022) gives a characterisation of what substructures Subgraph GNNs  
 112 cannot count at node-level based on the notation of cycles and paths. They show that the Subgraph  
 113 GNNs cannot count cycles of four or more nodes and paths of three or more nodes.

114 On a similar note, Huang et al. (2022) propose extracting edge-rooted subgraphs rather than node-  
 115 rooted and marking nodes that form the edge that anchors the subgraph extraction. Huang et al.  
 116 (2022) prove that the utilisation of double marking grants enhanced computational capabilities com-  
 117 pared to node rooted Subgraph GNNs. Additionally, it is ascertained that their model exhibit partial  
 118 superiority over the 3-WL test, enabling the counting of cycles with lengths shorter than seven and  
 119 all subgraphs up to size four in a non-induced setting.

120 Recently a new theoretical view of Subgraph GNNs based on the Subgraph Weisfeiler-Lehman,  
 121 a new version of the WL test, has been proposed (Zhang et al., 2023). This analysis presents a  
 122 characterisation of the expressive power of all node-rooted Subgraph GNNs. They conclude that  
 123 no node-rooted Subgraph GNN can be more powerful than the 2-folklore-WL (3-WL). This bound  
 124 was already discussed by Yan et al. (2023) and Frasca et al. (2022). However, Zhang et al. (2023)  
 125 demonstrate that no node-rooted Subgraph GNN can achieve the maximum expressivity of their time  
 126 complexity class. This result draws a limitation in the design of node-rooted Subgraph GNNs. Later  
 127 work by Yan et al. (2023) characterise the counting power of Subgraph GNNs for general architec-  
 128 tures and a general number of rooted nodes used as backbone. Furthermore, a method to compare  
 129 the expressivity of GNN models was introduced by Zhang et al. (2024) based on homomorphisms,  
 130 they summarise the ability of multiple GNN models on their ability to count any substructures with  
 131 no more than eight edges and no more than six vertices. Regarding induced subgraph counting at  
 132 graph level, the subject of our work, 1-WL models cannot count any pattern with 3 or more nodes.

### 133 3.3 STEP THREE - HOW IS SIGNIFICANCE OBTAINED

134  
 135 After obtaining the frequency of the structures in  $S_G$ , the next step is to evaluate their significance.  
 136 Hence, it is necessary to have an idea of what would produce, with no factor other than random  
 137 chance, a network similar to  $\mathcal{G}$  for some characteristic of interest. Let us denote NULL as a model  
 138 that can achieve that goal. One example of NULL is a model that, given a graph  $\mathcal{G}$ , randomly  
 139 switches edges while keeping the degree distribution of  $\mathcal{G}$  – degree distribution is the characteristic  
 140 of interest. The rewiring process is completely random and without any bias towards any predispo-  
 141 sition (Milo et al., 2004a;b).

### 142 3.4 MOTIFS AND GRAPH NEURAL NETWORKS

143  
 144 Motif estimation, when approached through the lens of GNNs, appears to be a challenge that, to the  
 145 best of our knowledge, remains largely unexplored in the existing literature.

146  
 147 **Directly counting.** One of the approaches that better matches direct motif estimation with GNNs  
 148 would be to count subgraphs of interest,  $S_G$ , using a GNN in the input graph  $\mathcal{G}$  and, after selecting  
 149 a suitable null model, generate an amount  $T$  of graphs according to it and use the same GNN model  
 150 used in  $\mathcal{G}$  to count the occurrences of the selected set of subgraphs in each of the  $T$  control graphs.

151 **Motifs as tool.** Other works that integrate GNNs and motifs typically deviate from estimating motifs  
 152 and use pre-computed ones to enhance the power of GNNs. Examples of this work include Motif  
 153 Convolutional Networks (Lee et al., 2018), motif2vec (Dareddy et al., 2019), Motif Graph Attention  
 154 Network (Sheng et al., 2024), Motif Graph Neural Network (Chen et al., 2023) and Heterogeneous  
 155 Motif Graph Neural Network (Yu & Gao, 2022). In the field of GNNs, another usage of the concept  
 156 of a motif as a relevant pattern comes from the attempt to explain the decision of GNN models. Two  
 157 examples in this field are GNNExplainer (Ying et al., 2019) and TempME (Chen & Ying, 2023).

#### 158 3.4.1 LEARNING RELEVANT PATTERNS

159  
 160 We will introduce the forthcoming studies under the term “relevant patterns” since most of them use  
 161 a definition of motif that is different from the one we introduce. Nonetheless, when discussing such  
 works, we follow the terminology of the original works and will call the pattern “motif”.

Works directly pertaining to motif estimation (used to refine a downstream task) are MICRO-Graph (Zhang et al., 2020) and MotiFiesta (Oliver et al., 2022). An example of a method made to estimate subgraph frequency is SPMiner (Ying et al., 2020).

The main problems of these works are the following: (1) either the model does not assume a null model and returns raw counts of occurrences of a general  $\mathcal{H}$  in  $\mathcal{G}$  (SPMiner and other frequency estimation models), or (2) the model may use a null model to guide motif search but only returns the subgraph(s) that are considered a motif by the model, meaning it is typically not possible to query for a specific  $\mathcal{H}$  (MICRO-Graph and MotiFiesta). Hence, these models typically ignore everything not branded as a motif, sometimes not even returning a motif score for the graphs regarded as such. Furthermore, models that return the raw count of occurrences can suffer from poor generalisation since the number of graph structures grows super-exponentially (Fu et al., 2023). As the size of a graph  $\mathcal{G}$  grows, the possible counts of a substructure  $\mathcal{H}$  in  $\mathcal{G}$  also grow super-exponentially, causing high variation between results of small and very-large graphs. This fact can hinder the learning process of models that aim at being agnostic of network size and topology. Also, for raw count models, since no null model is assumed, obtaining significance implies added subsequent computation.

## 4 INITIAL PROBLEM DESCRIPTION

Hereafter, referencing the number of occurrences of a graph  $\mathcal{H}$  within  $\mathcal{G}$ , denotes the induced count of  $\mathcal{H}$  in  $\mathcal{G}$ . Furthermore, all graphs are undirected and they do not have edge features.

According to the definition of motif adopted, to understand if a graph  $\mathcal{H}$  is a motif of a graph  $\mathcal{G}$ , we must know the number of occurrences of  $\mathcal{H}$  in  $\mathcal{G}$ . Let us denote such count as  $C(\mathcal{H}, \mathcal{G})$ . Furthermore, to grasp the importance of  $\mathcal{H}$  in  $\mathcal{G}$ , it is needed to know the count of  $\mathcal{H}$  across sufficient graphs derived from a null model denoted as NULL. Let us denote the average count of  $\mathcal{H}$  in graphs derived from NULL as  $C^\mu(\mathcal{H}, \mathcal{G}_{\text{NULL}})$  and the standard deviation as  $C^\sigma(\mathcal{H}, \mathcal{G}_{\text{NULL}})$ . Hence,  $Z(\mathcal{H}, \mathcal{G}_{\text{NULL}}) = \frac{C(\mathcal{H}, \mathcal{G}) - C^\mu(\mathcal{H}, \mathcal{G}_{\text{NULL}})}{C^\sigma(\mathcal{H}, \mathcal{G}_{\text{NULL}})}$  denotes the standardization (Z-score) of the occurrences of a graph  $\mathcal{H}$  in  $\mathcal{G}$ .

### 4.1 OUR APPROACH

Even though not used for subgraph counting, we implant degree features into the graphs to enhance the capability of the models. Instead of modelling the learning task as predicting a single value  $Z(\mathcal{H}, \mathcal{G})$  for some  $\mathcal{H}$  and some  $\mathcal{G}$ , we model it as a multi-target regression problem in order to predict the motif score of multiple subgraphs at once. This characterisation naturally allows the construction of motif fingerprints as proposed by Milo et al. (2004b). Thus, we start with a vector of Z-scores,  $z = [Z(\mathcal{H}_1, \mathcal{G}) \dots Z(\mathcal{H}_n, \mathcal{G})]$ .

The restriction of the number of graphs in  $z$  implies that the proposed model will not be able to search if an arbitrary graph is or is not a motif. However, by having a model that has a more restricted objective, we aim to achieve higher precision in the said objective. Since  $z$  has a restricted size, one other aspect that deserves careful consideration is deciding what is the size of  $z$  and what graphs compose it. Should the selected graphs exhibit negligible relation, an attempt to predict the Z-score concurrently for all graphs may prove harmful to the performance of the model. In this case, such an approach forces the model to incorporate distinct patterns to predict scores for each graph, thereby resulting in a sub-optimal global predictive efficacy. However, if  $z$  is composed of a well-thought group of graphs, allowing them to share common patterns from a learning perspective, we hypothesise that the performance of the model can improve when compared to predicting just one Z-score, due to the possibility of what is learned about a target variable be “shared” with others through weight sharing (one other advantage is the need to only train a single model instead of multiple). A good candidate for  $z$  should have patterns that are interconnected with each other, either from the point of view of the Z-score distribution or from a topological one.

Building on top of what was described in the two previous paragraphs, we focus on small graphs, in particular all connected graphs of size three and four. This is also supported by existing relevant literature (Milo et al., 2004b;a; Shen-Orr et al., 2002; Asikainen et al., 2020; Pržulj, 2007; Ribeiro & Silva, 2013) suggesting that to understand a complex network, it is an important to understand how small graphs behave in that network. We focus on these graphs because their proximity in size should allow them to have a topological connection that translates in a connection in their Z-Scores.

Restricting the size of the graphs used in  $z$  to small ones also has the added benefit that we can get the ground-truth of motif scores for a diverse type and size of networks, allowing for a richer train dataset. Furthermore, we expect that using a set of graphs of increasing size in the number of nodes and edges gives enough interconnectivity between their patterns from both a topological and a Z-score distribution point of view to allow the model to have a strong inductive bias towards meaningful patterns, allowing for a stronger performance. For example, a graph with many size four cliques will probably have a small amount of 4-stars.

We will refer to the set of graphs used to populate  $z$  as  $\Omega$ . The function notation  $\Omega(\mathcal{G})$  gives the set of all graphs that have the same number of connected nodes as the graph  $\mathcal{G}$ . For the chosen graphs it is possible to create two groups in  $\Omega$ , the graphs of size three and the graphs of size four. Using the motif Z-score directly in the learning task not only makes the function and result highly interpretable, but also eliminates the need to compute multiple networks based on the NULL model to determine significance. By normalising the Z-score across groups of graphs, as defined in Milo et al. (2004b), the values of  $z$  are constrained between  $-1$  and  $1$ , independent of network size, enhancing the model’s predictive stability across various network scales. Moreover, normalising the Z-scores with  $s_i = z_i / (\sum_{j \in \Omega(i)} z_j^2)^{1/2}$  imposes a mathematical interconnectivity between the Z-scores of graphs of the same group. This relationship, where the sum of squared Z-scores equals 1, supports a multi-target objective and further strengthens the problem formulation by adding an additional layer of interdependence among graphs. The normalised Z-score refers to the values of the significance-profile  $s$ . The learning task consists of minimising the MSE between the true and predicted significance-profiles.

#### 4.2 ON THE RELATION WITH EXPRESSIVITY THROUGH SUBGRAPH COUNTING

It is expected that the expressivity regarding substructure counting to be highly related to the expressivity of discovering the significance-profile of graphs. Concretely, the problem of counting graphs is a subset of the problem of discovering significance-profiles where reducing the null model to nothing reduces the problem to graph counting.

Since  $P$ , the problem of counting graphs, is a subset of  $S$ , the problem of significance-profile estimation, it is possible to obtain instances of  $S$  that are as hard as  $P$ , easier than  $P$  and harder than  $P$ . Under the assumption that  $S$  and  $P$  function around the same set of graphs, these differences in difficulty come from the choice of null model. In the case of  $S = P$ , the null model should do nothing, for example, returning always 0. In the case of  $S < P$ , the null model could always return the counts of each subgraph in a graph  $\mathcal{G}$  without modifying  $\mathcal{G}$ , reducing the problem to always predicting a vector of zeros. For the case of  $S > P$ , employing a null model that randomly returns counts for  $\mathcal{G}$  should make the problem theoretically harder since the model would have to learn the random process employed to correctly construct the significance-profiles. Thus, theoretical guarantees of expressivity might not hold depending on the selected null model. For instance, a recent demonstration solved the dimensionality of the  $k$ -WL test for induced subgraph count, stating that to perform induced subgraph count of any pattern with  $k$  nodes we need at least dimensionality  $k$  (Lanzinger & Barceló, 2023). Furthermore, no induced pattern with  $k + 1$  nodes can be counted with dimensionality  $k$ , a result not verified to non-induced counts (Lanzinger & Barceló, 2023). However, depending on the null model, when working with significance-profiles over graphs of size four in  $\Omega$ , it might not be enough to have a model as powerful as the 4-WL to guarantee enough power to correctly identify the normalised Z-scores of the size four graphs in  $\Omega$ .

**Modifications made to the problem.** In Section 4, we attempted to reduce the difficulty of the problem by formulating it as a multitarget regression of interconnected values from an algebraic and topological point of view. However, even in the case where it exists a perfect dependence between graphs of the same group of  $\Omega$ , or even across different groups, the problem is still at least as hard as finding the significance-profile of the graph(s) that governs the dependence relation. For example, if the significance-profile of the 3-path was symmetric to the significance-profile of the 3-clique, it would still be necessary to determine the significance-profile for 3-paths to compute the significance-profile of 3-cliques and vice versa.

**Testing with 1-WL bounded models?** MPNNs cannot perform induced counts of patterns of three or more nodes (Chen et al., 2020). Nevertheless, MPNNs are not inherently incapable of counting patterns in any graphs. Rather, for a pattern  $\mathcal{H}$ , there exists graphs  $\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{G}$  such that  $C(\mathcal{H}, \mathcal{G}_1) \neq$

270  $C(\mathcal{H}, \mathcal{G}_2)$  and for any MPNN  $M$  under 1-WL,  $M(\mathcal{G}_1) = M(\mathcal{G}_2)$ . Hence,  $M$  cannot discover  
 271  $C(\mathcal{H}, \mathcal{G}_1)$  and  $C(\mathcal{H}, \mathcal{G}_2)$  simultaneously. However, within the 1-WL framework, MPNNs remain  
 272 highly valuable and find numerous practical applications in real-world scenarios. Furthermore, we  
 273 did not construct any characterisation of the problem space of  $S$  regarding  $P$  for the null model used.  
 274 Hence, we might have made the problem easier (or harder) than substructure counting. The same  
 275 applies to the usage of a multi-target objective. Thus, we will scrutinise the capability of MPNNs to  
 276 address our particular challenge. Furthermore, more expressive models like Subgraph GNNs have a  
 277 very high computational complexity, being very hard to use in large-scale graphs.

## 278 5 DATASETS

279 We rely exclusively on synthetic graphs, rather than commonly used GNN datasets such as Wu et al.  
 280 (2018); Gómez-Bombarelli et al. (2018); Morris et al. (2020); Hu et al. (2021b;a). Our goal is to  
 281 have a large dataset with diverse topological features, which would be difficult and costly to obtain  
 282 from real-world data spanning multiple domains. Using multiple graph generators designed to sim-  
 283 ulate real-world phenomena and fully exploring their parameter space, we aim to create a dataset  
 284 with both high topological diversity and a close resemblance to real-world data. Additionally, syn-  
 285 thetic data allows for flexible expansion of the dataset size. Using synthetic data to train GNNs is  
 286 not a new concept. However, as far as we know, most of the popularly used datasets typically have  
 287 very small graphs (at most few hundreds of nodes) and are generated from a small set of generators,  
 288 often random regular graphs and Erdős-Renyi graphs (Chen et al., 2020). Another popular type of  
 289 synthetic graph dataset for benchmarking is small handcrafted graphs to limit test GNN models (Ab-  
 290 boud et al., 2021; Murphy et al., 2019; Balcilar et al., 2021; Wang & Zhang, 2023). While still very  
 291 limited, the only exceptions identified are Veličković et al. (2020); Corso et al. (2020).

292 Since we are not interested in limit testing the power of the model in comparison to theoretical  
 293 tests, and instead aim at having a diverse dataset regarding graph topologies and motif scores, we  
 294 create a new dataset using a total of 23 synthetic generators (11 non-deterministic, 12 deterministic).  
 295 We explore their graph-generating space in order to extract all types of possible topologies they  
 296 can express while limiting the graph size in order to avoid bottlenecks that increase training time  
 297 beyond what we find reasonable for the amount of time and resources available. The final dataset  
 298 puts the non-deterministic segment with 109164 graphs, and the deterministic segment with 38400  
 299 graphs, totalling for  $\approx 250$  million nodes and  $\approx 750$  million edges. Section A has a description  
 300 with greater detail on how each generator was explored. For the ground-truth, we calculate  $s$  using  
 301 G-Tries (Ribeiro & Silva, 2013).  
 302

303 Exploring the myriad of significance-profiles from the generated graphs leads us to the con-  
 304 clusion that the 3-path and the triangle can only take on a few sets of values, being those  
 305  $\{-0.707106, 0, 0.707106, 1\}$  for the 3-path and the first three values for the triangle. This leads  
 306 to a strict result on the inter-variability of their Z-scores. Apart from cases where both scores are  
 307 zero or the 3-path is 1 and the triangle is 0 (an artefact from the G-Trie model, primarily affecting the  
 308 duplication-divergence model), we found that the Z-scores for graphs of size three are symmetric.  
 309 This means that if one structure takes a Z-score of  $x$ , the other will take  $-x$ . In normalised form,  
 310 these values are mapped to 0.707106 and  $-0.707106$ , respectively. For the size four graphs, we  
 311 can say that the size three encodes some information about the significance-profile of graphs of size  
 312 four, alluding to the possibility of an advantage of using graphs of both sizes in the target variable  
 313 (Figure 3 - Appendix A.4 has a detailed view of this result). As for strict dependence between the  
 314 graphs of size four, apart from the mathematical based constraint, we could not confirm any other.

315 **Real-World Data.** Since we are interested in assessing the performance of the models with real-  
 316 world data, we compiled a test set based on real networks of multiple categories. Besides varying  
 317 the type of network, we vary in their relative size in terms of number of nodes and edges. We devise  
 318 two categories based on the size of the networks used in the train set: (1) small-scale networks, (2)  
 319 medium-large scale networks. Section A.5 presents a detailed description of the networks collected.

## 320 6 METHODOLOGY

321 The model used in the experiments is a very simple model similar to the one described in Chen et al.  
 322 (2020) definition A.1. from Appendix A. Section B further details the model and how it was trained.  
 323

**Baselines.** We define a baseline denoted by an horizontal red line corresponding to the expected loss incurred when predicting the significance-profile  $s$  by employing an independent random uniform model for each component of  $s$ . A second base line denoted by an horizontal blue line corresponds to the expected loss when using a model predicting a random value for every component of  $s$ , but taking into account the restrictions posited in Section 4 regarding the range of values each group of  $\Omega$  can take. A final baseline dependent on the architecture of the model is represented by a light pink bar. It illustrates the mean error derived from using the model in question with random weights.

**Persistent Patterns.** Given a set of significance-profiles  $\mathbb{S}$ , we define the persistency  $\rho$  of a pattern  $s \in \mathbb{S}$  as the frequency of its occurrence within  $\mathbb{S}$ . The higher the number of occurrences of  $s$ , the more persistent the pattern will be. Considering the randomness inherent to deriving a significance-profile, we use a threshold,  $t$ , to decide what patterns are equivalent to each other. Given  $t$ , discovering the persistence of  $s$  is reduced to  $\rho_s = |\{s' \mid d_\infty(s, s') \leq t, \forall s' \in \mathbb{S}\}|$  where  $d_\infty$  is the Chebyshev distance. We employ the concept of persistent patterns to develop a high-level understanding of the nature of the true and predicted significance-profiles. We use the true significance-profiles as a standard to determine an appropriate threshold for equivalency. This threshold is then applied to the predicted significance-profiles. Section C further details the process to obtain  $t$ .

**“Correct” Predictions.** Since having a criteria for a prediction for a significance-profile to be correct/useful depends largely on the research field, in order to have a systematic threshold, we study how many predictions for each generator have all individual values of the predicted significance-profile below the obtained test error. Furthermore, we ensure that for a prediction to count as “correct”, all significance-profiles for graphs in  $\Omega$  have the correct sign. Presuming the test errors obtained for the synthetic dataset are relatively low, this approach guarantees that the predictions deemed “correct” not only contribute to a lower test loss than the one observed but are also potentially valuable from an application perspective. Let us denote a benchmark model as  $T$ . Let this model predict a random pattern according to the criterion used for the horizontal blue line described above. If we define a cutoff mean loss for  $T$  of 0.25 (0.5 absolute error), this model will incur in a possible error of 25% of the maximum allowed. According to  $1e8$  simulations,  $T$  will have a rate of “correct” guesses of 0.364% while following the restrictions here introduced.

## 7 RESULTS

The chosen models correspond to one instance of GIN and SAGE for the non-deterministic segment and one instance of GIN for the deterministic one (see Appendix D for more details). Figure 1 shows the results for the selected models for the small and medium-large datasets. The yellow bar denotes the validation error, the blue bar represents the test error observed in the test dataset and the green bar the result of the model in the real-world dataset. The other marks represent the baselines from Section 6. A more detailed version of the predictions is available in Appendix D.2.

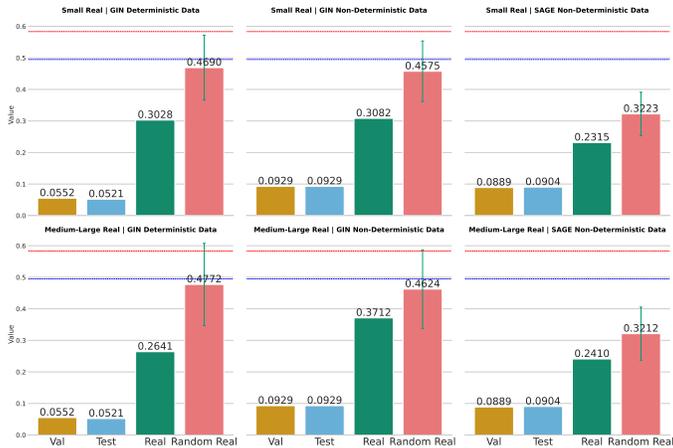


Figure 1: Scores of the best models for the test and real-world datasets.

## 7.1 DISCUSSION OF THE RESULTS

The models have low generalisation error as training, validation, and test errors are similar, and there is no significant difference in performance between GIN and SAGE for synthetic data. However, the models present a much higher error in the real-world dataset than the out-of-sample error estimated through the test set. Furthermore, the obtained value is dangerously close to the error obtained with random weights minus one standard deviation. Nonetheless, one positive is the stability of the error across the two different real-world dataset sizes. Even though the error is high, it remains stable as the size of the networks increases, a growth that exceeds three orders of magnitude in scale.

**What is the model not doing.** It is evident that any trained model does not adhere to an independent uniform random prediction and is not randomly predicting based on the restrictions of  $s$ . All predictions conspicuously surpass the red line and blue line, established as the initial benchmark.

**What can the model be doing.** Regarding the poor performance, one possible scenario is that the synthetic data used does not accurately reflect the real-world data, leading to the observed dissonance between scores in the test and real datasets. However, it is essential to consider the limitations of the MSE for this particular analysis. While MSE is valuable for comparing models with significant differences in performance, it does not provide sufficient information regarding individual predictions and their global shape. For instance, a model predicting a value of 0 for all graphs would have an expected MSE of 0.25. Using the concept of persistent patterns, we can understand how the predictions are distributed in terms of proximity to one another in relation to what would be expected from the ground-truth. Table 3 (Appendix C) conveys a detailed view of the clusters found for the threshold inferred from the true significance-profiles. Summarily, we conclude: (1) all models exhibit a substantially lower number of predicted patterns than expected; (2) this reduction is primarily attributed to the higher  $\rho$  of select patterns, rather than an higher mean  $\rho$  across patterns.

In the following sections we answer: (1) What can cause a model to tend towards persistent patterns? (2) Can we take the low error in the test set as a signal that the model is learning the synthetic data? (3) What causes the discrepancy between the scores in the synthetic and the real-world dataset?

## 7.2 TENDENCY FOR PERSISTENT PATTERNS.

This study examines the impact of dropout on the expressivity of models and persistence of the patterns. By systematically increasing dropout values, the model’s power is reduced due to regularisation. Testing with both GIN and SAGE models reveals that dropout on the MPNN has little to no effect on GIN’s performance but slightly affects SAGE, indicating GIN’s greater expressivity. This possibility was briefly discussed in Section 4.2, suggesting that less expressive GNN models are capable of distinguishing fewer graphs. The results suggest that models lacking expressivity generate fewer diverse patterns, leading to similar predictions across graphs. Hence, poor performance on real-world data might not solely result from limitations of the synthetic dataset, but also from the lack of expressivity of the models.

**Answering question (1).** Tendency towards persistency of patterns higher than what would be expected stems from the lack of expressivity of either the model, the dataset, or both.

## 7.3 MODEL PREDICTIONS IN THE SYNTHETIC DATASET

Following Figure 8 (Appendix D.2), we conclude that the predictions made by all the models are reasonable for all generators. The results suggest that the model is sufficiently expressive to distinguish between different graph generators, as predictions often align with the true mean significance-profile. However, it struggles to differentiate graphs within the same generator, particularly in non-deterministic and highly diverse generators like random regular and Erdős-Renyi. The differences between graphs of different generators is large enough that an MPNN can capture them. However, as we see the predictions gravitating towards the mean pattern of a generator, evidenced by the tighter band between the 2.5% and 97.5% percentile, distinguishing between graphs within each generator seems to be a task that MPNNs cannot generally perform. In these cases, the model tends to predict a less diverse set of patterns, leading to tighter percentile bands compared to reality (result also observable through Table 3 in Appendix C).

Table 1: Number of graphs with significance-profiles deemed “correct”. *Italic-underlined* are for SAGE and the others GIN. Generators with  $\geq 50\%$ ,  $\geq 70\%$  and  $\geq 90\%$  are highlighted. Error in non-deterministic graphs for SAGE is 15.0%, 15.2% for GIN and 11.4% for deterministic ones.

Generator	“Correct”		“Incorrect”		Generator	“Correct”	“Incorrect”
Geometric-3D DD Graph	<i>035</i>	<i>224</i>	<i>314</i>	125	Balanced Tree	<i>191</i>	129
Duplication Divergence Graph	<i>280</i>	<i>303</i>	<i>069</i>	046	Barbell Graph	<i>185</i>	135
Extended Barabasi Albert Graph	<i>152</i>	056	<i>197</i>	293	Binomial Tree	<i>197</i>	123
Erdős-Renyi	<i>000</i>	000	<i>349</i>	349	Chordal Cycle Graph	076	244
Forest Fire	<i>209</i>	119	<i>140</i>	230	Circular Ladder Graph	158	162
Gaussian Random Partition Graph	<i>000</i>	003	<i>349</i>	346	Dorogovtsev Goltsev Mendes Graph	107	213
Random Limited Geometric Graph	<i>006</i>	085	<i>343</i>	264	Full Rary Tree	<i>244</i>	076
Newman Watts Strogatz Graph	<i>000</i>	<i>189</i>	<i>349</i>	160	Square Lattice	<i>167</i>	153
Powerlaw Cluster Graph	<i>336</i>	<i>301</i>	<i>013</i>	048	Hexagonal Lattice Graph	156	164
Random Regular Graph	<i>000</i>	000	<i>349</i>	349	Lollipop Graph	085	235
Watts Strogatz Graph	<i>000</i>	031	<i>349</i>	318	Star Graph	<i>106</i>	052
					Triangular Lattice Graph	<i>260</i>	060
Total	<i>1018</i>	1311	<i>2821</i>	2528	Total	1932	1746

Assuming that the choice of null model had little impact in the difficulty of the problem, the conclusion of the ability of a model with expressivity at most equal to the 1-WL to be able to distinguish graphs of different generators can be seen as a partial empirical confirmation of an old result by Babai & Kucera (1979); Babai et al. (1980), regarding the 1-WL test being able to distinguish any random graph with high probability as the size of graph approaches infinity. Similarly, the apparent good performance of the tree generators is also theoretically backed by findings in Arvind et al. (2015). The conclusion of the inability of the model to distinguish graphs with high granularity among those in the same generator has theoretical backing for the case of the random regular generator (Babai & Kucera, 1979; Cai et al., 1989; Babai et al., 1980). As for the other generators, following the result in Babai & Kucera (1979); Babai et al. (1980), theoretically, it should be highly probable that a model as powerful as the 1-WL could distinguish most of the graphs, specially random ones, at inter and intra-generator level. However, according to our findings, this is not exactly true in practice, either because the model could not reach the 1-WL expressivity due to failing to approximate an injective function or because the bound does not work well in practical scenarios.

The model may be capable of accurately making inter-generator predictions by predicting a significance-profile corresponding to the mean of each generator, along with additional predictions gravitating around it. Following the formulation from Section 6, the “correct” predictions using the test set error will be evaluated assuming an error between 15.4% and 11.4% of the total error.

According to Table 1, more than 50% of the predictions for the deterministic dataset are satisfactory/correct. Not counting the regular graphs, known to not be distinguishable by 1-WL, SAGE got satisfactory predictions for 29.2% of the graphs and GIN for 37.6%. The generators whose graphs had more satisfactory predictions were those exhibiting a stronger mean pattern. In these cases, a model reaps significant gains from simply following the mean pattern. For instance, if we do not count the three problematic generators, regular graph, Gaussian random partition, and Erdős-Renyi, GIN has a satisfactory prediction rate of 46.9% and SAGE 36.5%. The worst-performing model achieves 29.2%, more than 80 times better than  $T$ , even with its lenient margin for accuracy.

**Answering question (2).** Yes, the model is learning the synthetic dataset (further reinforced by Figure 8). The model learns in two scales: inter-generator and intra-generator. It presents a good general capacity for inter-generator learning, meaning it does a fine job of distinguishing what is the generator of a graph. As for its intra-generator performance, it has a reasonable discriminative power. In the best case, for the non-deterministic segment, assuming 15.4% margin of error and a guaranteed signal match, it guesses correctly the significance-profile 46.9% of the times. As for the deterministic segment it predicts correctly 52.6% of the times, with a margin of error of 11.4%.

#### 7.4 MODEL PREDICTIONS IN THE REAL-WORLD DATASET

Analysing the multiple figures from Section D.2, it becomes apparent that the inter-category (or inter-generator, in the context of the synthetic dataset) performance is far from optimal. The models exhibit identical predictions across graphs from different categories, even when these categories have distinct significance-profiles. Intra-category performance is even poorer than inter-category, with models like SAGE generating patterns that are too similar across graphs, though just distinct enough to avoid being captured by the persistence measures. This phenomenon is exemplified by

the predictions for the small biological category and the small interaction category observed in the SAGE model, by the medium-large interaction and medium-large social communication for GIN non-deterministic, among others. In the real-world data analysis, the SAGE model had the best performance with 12.5% satisfactory predictions, seven correct in the small dataset. Other models had three or fewer. This result is  $2.4\times$  worse than the poorest in the synthetic dataset, despite a 24.1% margin of error,  $\approx 1.6\times$  larger.

Interestingly, the models seem to have induced their own groupings based on similarities with the synthetic datasets used for training. For instance, real-world graphs like *ia-escorts-dynamic*, *coauthor-CS*, and *ia-primary-school-proximity* exhibited patterns resembling synthetic models like duplication-divergence, forest-fire, and geometric graphs, respectively. This suggests that the model predicts based on how similar a real-world graph is to the synthetic ones it has seen. While the model struggles to produce satisfactory predictions for real-world networks, it could help identify which synthetic model a real-world network most closely resembles. The closer the predicted pattern is to the true pattern, the more alike the synthetic and real networks.

**Points of divergence.** In network similarity discovery based on significance-profiles, two key concepts are crucial. First, the model’s ability to distinguish networks is constrained by the expressivity of the space of the significance-profile used. The more expressive the space, the more reliable the ability of the model to distinguish networks. Secondly, if the model predicts similar profiles for two graphs  $\mathcal{G}$  and  $\mathcal{H}$ , indicating they resemble a graph  $\mathcal{F}$ , this suggests  $\mathcal{G}$  and  $\mathcal{H}$  may originate from a process similar to  $\mathcal{F}$ . However, this conclusion is valid only if the true profiles of  $\mathcal{G}$  and  $\mathcal{H}$  are indeed similar; otherwise, the model’s lack of expressivity leads to incorrect conclusions.

**Closeness to Random Weights.** A network with random weights can take any possible network as a value. However, achieving performance near the score of the trained network purely by chance is unlikely due to the complexity of the model. To address this, we propose that while the solution space for random models may differ, its mean score should be similar to that of the solution space of the trained model. Specifically, (1) models with random weights resemble those with high dropout, showing a tendency towards highly persistent patterns; (2) the trained model produces meaningful patterns for synthetic data; (3) the predictions of the model can significantly differ from true significance-profiles due to its limited expressivity. Thus, while a trained model predicts meaningful patterns ineffectively, a random one only predicts highly persistent patterns, leading to a similar average scores.

**Answering question (3).** The model learns from the synthetic dataset. The discrepancy comes from the synthetic data not accurately reflecting the real-world, at least when used by models limited by the 1-WL. Thus, we can only confirm that the knowledge extracted from the synthetic dataset by the used models is not enough to describe real-world data.

## 7.5 VALIDATION OF THE ASSUMPTIONS MADE

Multi-target regression improves predictive accuracy for most graphs, except for graphs 4-clique and 4-path, likely due to limited benefit from shared information. Predicting significance profiles directly enhances performance for graphs with significant variation in subgraph occurrences, like the 3-path, triangle, 4-path, and 4-clique. Despite minor error increases in some cases, the overall gains in accuracy and computational efficiency make this approach preferable for motif estimation in the given null model. Section D.1 has a detailed exploration of the mentioned results.

## 8 CONCLUSIONS

Despite the lack of a GNN-based method specifically designed for predicting significance profiles, our current MPNN models combined with synthetic data are still deemed insufficient for real-world applications related to significance-profile discovery. The best performance benchmark shows a prediction accuracy of 12.5% for small networks and 10.41% for medium-large networks, each with a 24.1% margin of error. In contrast, having into account their simplicity, the used models together with the presented formulation, for synthetic data, achieve good results, with a benchmark accuracy of 46.9% with a 15.4% margin of error for non-deterministic generators and 52.6% with an 11.4% margin of error for deterministic generators. This suggests that the models are promising for network categorisation by effectively distinguishing high-level differences between graphs.

## REFERENCES

- 540  
541  
542 Ralph Abboud, İsmail İlkan Ceylan, Martin Grohe, and Thomas Lukasiewicz. The Surprising Power  
543 of Graph Neural Networks with Random Node Initialization. In *Proceedings of the Thirtieth*  
544 *International Joint Conference on Artificial Intelligence*, pp. 2112–2118, California, aug 2021.  
545 International Joint Conferences on Artificial Intelligence Organization. ISBN 978-0-9992411-9-  
546 6. doi: 10.24963/ijcai.2021/291.
- 547 Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna:  
548 A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM*  
549 *SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.  
550
- 551 Réka Albert and Albert-László Barabási. Topology of evolving networks: Local events and uni-  
552 versality. *Physical Review Letters*, 85(24):5234–5237, 2000. ISSN 10797114. doi: 10.1103/  
553 PhysRevLett.85.5234.
- 554 Réka Albert, Hawoong Jeong, and Albert-László Barabási. Diameter of the world-wide web. *Nature*,  
555 401(6749):130–131, September 1999. ISSN 1476-4687. doi: 10.1038/43601.
- 556  
557 Uri Alon. Network motifs: theory and experimental approaches. *Nature Reviews Genetics*, 8(6):  
558 450–461, jun 2007. ISSN 1471-0056. doi: 10.1038/nrg2102.
- 559  
560 V. Arvind, Johannes Köbler, Gaurav Rattan, and Oleg Verbitsky. *On the Power of Color Refinement*,  
561 pp. 339–350. Springer International Publishing, 2015. ISBN 9783319221779. doi: 10.1007/  
562 978-3-319-22177-9\\_26.
- 563 Aili Asikainen, Gerardo Iñiguez, Javier Ureña-Carrión, Kimmo Kaski, and Mikko Kivelä. Cumu-  
564 lative effects of triadic closure and homophily in social networks. *Science Advances*, 6(19), May  
565 2020. ISSN 2375-2548. doi: 10.1126/sciadv.aax7310.
- 566  
567 Laszlo Babai and Ludik Kucera. Canonical Labelling of Graphs in Linear Average Time. *An-*  
568 *annual Symposium on Foundations of Computer Science - Proceedings*, (2):39–46, 1979. ISSN  
569 02725428. doi: 10.1109/sfcs.1979.8.
- 570  
571 László Babai, Paul Erdős, and Stanley M Selkow. Random Graph Isomorphism. *SIAM Journal on*  
572 *Computing*, 9(3):628–635, aug 1980. ISSN 0097-5397. doi: 10.1137/0209047.
- 573 Muhammet Balcilar, Pierre Heroux, Benoit Gauzere, Pascal Vasseur, Sebastien Adam, and Paul  
574 Honeine. Breaking the limits of message passing graph neural networks. In Marina Meila and  
575 Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, vol-  
576 ume 139 of *Proceedings of Machine Learning Research*, pp. 599–608. PMLR, 18–24 Jul 2021.
- 577  
578 Albert-László Barabási and Márton Pósfai. *Network science*. Cambridge University Press, 2017.
- 579  
580 Albert-László Barabási, Réka Albert, and Hawoong Jeong. Scale-free characteristics of random  
581 networks: the topology of the world-wide web. *Physica A: Statistical Mechanics and its Applica-*  
582 *tions*, 281(1-4):69–77, June 2000. ISSN 0378-4371. doi: 10.1016/s0378-4371(00)00018-2.
- 583 Alexis Bénichou, Jean-Baptiste Masson, and Christian L. Vestergaard. Compression-based inference  
584 of network motif sets. pp. 1–51, nov 2023.
- 585  
586 Ulrik Brandes, Marco Gaertler, and Dorothea Wagner. *Experiments on Graph Clustering Al-*  
587 *gorithms*, pp. 568–579. Springer Berlin Heidelberg, 2003. ISBN 9783540396581. doi:  
588 10.1007/978-3-540-39658-1\\_52.
- 589  
590 J.-Y. Cai, Martin Furer, and Neil Immerman. An optimal lower bound on the number of variables  
591 for graph identification. In *30th Annual Symposium on Foundations of Computer Science*, pp.  
592 612–617. IEEE, 1989. ISBN 0-8186-1982-1. doi: 10.1109/SFCS.1989.63543.
- 593  
Jialin Chen and Rex Ying. TempME: Towards the Explainability of Temporal Graph Neural Net-  
works via Motif Discovery. (NeurIPS):1–24, oct 2023.

- 594 Kaixuan Chen, Shunyu Liu, Tongtian Zhu, Ji Qiao, Yun Su, Yingjie Tian, Tongya Zheng, Haofei  
595 Zhang, Zunlei Feng, Jingwen Ye, and Mingli Song. Improving Expressivity of GNNs with  
596 Subgraph-specific Factor Embedded Normalization. In *Proceedings of the 29th ACM SIGKDD*  
597 *Conference on Knowledge Discovery and Data Mining*, volume 1, pp. 237–249, New York, NY,  
598 USA, aug 2023. ACM. ISBN 9798400701030. doi: 10.1145/3580305.3599388.
- 599 Zhengdao Chen, Lei Chen, Soledad Villar, and Joan Bruna. Can graph neural networks count  
600 substructures? *Advances in Neural Information Processing Systems*, 2020-Decem(NeurIPS),  
601 feb 2020. ISSN 10495258.
- 602 Stephen A. Cook. The complexity of theorem-proving procedures. In *Proceedings of the third*  
603 *annual ACM symposium on Theory of computing - STOC 71*, STOC 71. ACM Press, 1971. doi:  
604 10.1145/800157.805047.
- 605 Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal  
606 Neighbourhood Aggregation for Graph Nets. *Advances in Neural Information Processing Sys-*  
607 *tems*, 2020-Decem(NeurIPS), apr 2020. ISSN 10495258.
- 608 Jesper Dall and Michael Christensen. Random geometric graphs. *Physical Review E*, 66(1):016121,  
609 jul 2002. ISSN 1063-651X. doi: 10.1103/PhysRevE.66.016121.
- 610 Manoj Reddy Dareddy, Mahashweta Das, and Hao Yang. motif2vec: Motif Aware Node Rep-  
611 resentation Learning for Heterogeneous Networks. In *2019 IEEE International Conference*  
612 *on Big Data (Big Data)*, pp. 1052–1059. IEEE, dec 2019. ISBN 978-1-7281-0858-2. doi:  
613 10.1109/BigData47090.2019.9005670.
- 614 S. N. Dorogovtsev, A. V. Goltsev, and J. F. F. Mendes. Pseudofractal scale-free web. *Physical*  
615 *Review E*, 65(6), June 2002. ISSN 1095-3787. doi: 10.1103/physreve.65.066122.
- 616 Paul Erdős and Alfréd Rényi. On the Evolution of Random Graphs. *Magyar Tudományos Akadémia*  
617 *Értesítője*, 5(1):17–61, 1960.
- 618 Fabrizio Frasca, Beatrice Bevilacqua, Michael M. Bronstein, and Haggai Maron. Understanding and  
619 Extending Subgraph GNNs by Rethinking Their Symmetries. *Advances in Neural Information*  
620 *Processing Systems*, 35(NeurIPS), jun 2022. ISSN 10495258.
- 621 Tianyu Fu, Chiyue Wei, Yu Wang, and Rex Ying. DeSCo: Towards Generalizable and Scalable  
622 Deep Subgraph Counting. aug 2023.
- 623 A.A. Giannopoulos and V.D. Milman. Concentration property on probability spaces. *Advances in*  
624 *Mathematics*, 156(1):77–106, December 2000. ISSN 0001-8708. doi: 10.1006/aima.2000.1949.
- 625 Reid Ginoza and Andrew Mugler. Network motifs come in sets: Correlations in the randomization  
626 process. *Physical Review E*, 82(1), July 2010. ISSN 1550-2376. doi: 10.1103/physreve.82.  
627 011921.
- 628 Daniel Golovin, Benjamin Solnik, Subhodeep Moitra, Greg Kochanski, John Elliot Karro, and  
629 D. Sculley (eds.). *Google Vizier: A Service for Black-Box Optimization*, 2017.
- 630 Rafael Gómez-Bombarelli, Jennifer N. Wei, David Duvenaud, José Miguel Hernández-Lobato,  
631 Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D. Hirzel,  
632 Ryan P. Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven con-  
633 tinuous representation of molecules. *ACS Central Science*, 4(2):268–276, January 2018. ISSN  
634 2374-7951. doi: 10.1021/acscentsci.7b00572.
- 635 William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large  
636 graphs. In I Guyon, U Von Luxburg, S Bengio, H Wallach, R Fergus, S Vishwanathan, and  
637 R Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 2017-Decem, pp.  
638 1025–1035. Curran Associates, Inc., 2017.
- 639 Desmond J. Higham, Marija Rašajski, and Nataša Pržulj. Fitting a geometric graph to a protein-  
640 protein interaction network. *Bioinformatics*, 24(8):1093–1099, March 2008. ISSN 1367-4803.  
641 doi: 10.1093/bioinformatics/btn079.

- 648 Petter Holme and Beom Jun Kim. Growing scale-free networks with tunable clustering. *Physical*  
649 *Review E*, 65(2):026107, jan 2002. ISSN 1063-651X. doi: 10.1103/PhysRevE.65.026107.  
650
- 651 Weihua Hu, Matthias Fey, Hongyu Ren, Maho Nakata, Yuxiao Dong, and Jure Leskovec. Ogb-lsc:  
652 A large-scale challenge for machine learning on graphs, 2021a.
- 653 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta,  
654 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs, 2021b.  
655
- 656 Yinan Huang, Xingang Peng, Jianzhu Ma, and Muhan Zhang. Boosting the Cycle Counting Power  
657 of Graph Neural Networks with  $\mathbb{I}^2$ -GNNs. pp. 1–27, oct 2022.
- 658 I. Ispolatov, P. L. Krapivsky, and A. Yuryev. Duplication-divergence model of protein interaction  
659 network. *Physical Review E*, 71(6):061911, jun 2005. ISSN 1539-3755. doi: 10.1103/PhysRevE.  
660 71.061911.
- 661 Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Net-  
662 works. *5th International Conference on Learning Representations, ICLR 2017 - Conference Track*  
663 *Proceedings*, pp. 1–14, sep 2017.
- 664 Matthias Lanzinger and Pablo Barceló. On the power of the weisfeiler-leman test for graph motif  
665 parameters, 2023.
- 666 John Boaz Lee, Ryan A. Rossi, Xiangnan Kong, Sungchul Kim, Eunye Koh, and Anup Rao.  
667 Higher-order graph convolutional networks, 2018.
- 668 Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graph evolution: Densification and shrinking  
669 diameters. *ACM Transactions on Knowledge Discovery from Data*, 1(1), 2007. ISSN 15564681.  
670 doi: 10.1145/1217299.1217301.
- 671 Richard Liaw, Eric Liang, Robert Nishihara, Philipp Moritz, Joseph E Gonzalez, and Ion Sto-  
672 ica. Tune: A research platform for distributed model selection and training. *arXiv preprint*  
673 *arXiv:1807.05118*, 2018.
- 674 Alexander Lubotzky. *Discrete Groups, Expanding Graphs and Invariant Measures*. Birkhäuser  
675 Basel, 1994. ISBN 9783034603324. doi: 10.1007/978-3-0346-0332-4.
- 676 S. Mangan and U. Alon. Structure and function of the feed-forward loop network motif. *Proceedings*  
677 *of the National Academy of Sciences*, 100(21):11980–11985, October 2003. ISSN 1091-6490.  
678 doi: 10.1073/pnas.2133841100.
- 679 R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon. Network Motifs:  
680 Simple Building Blocks of Complex Networks. *Science*, 9781400841(October):217–220, 2004a.  
681 doi: 10.1515/9781400841356.217.
- 682 Ron Milo, Shalev Itzkovitz, Nadav Kashtan, Reuven Levitt, Shai Shen-Orr, Inbal Ayzenshtat, Michal  
683 Sheffer, and Uri Alon. Superfamilies of Evolved and Designed Networks. *Science*, 303(March):  
684 1538–1542, 2004b. doi: 10.1126/science.1089167.
- 685 Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang,  
686 Melih Elibol, Zongheng Yang, William Paul, Michael I. Jordan, and Ion Stoica. Ray: A distributed  
687 framework for emerging ai applications, 2018.
- 688 Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion  
689 Neumann. Tudataset: A collection of benchmark datasets for learning with graphs, 2020.
- 690 Ryan L. Murphy, Balasubramaniam Srinivasan, Vinayak Rao, and Bruno Ribeiro. Relational Pool-  
691 ing for Graph Representations. *36th International Conference on Machine Learning, ICML 2019,*  
692 *2019-June:8192–8202*, mar 2019.
- 693 M.E.J. Newman and D.J. Watts. Renormalization group analysis of the small-world network model.  
694 *Physics Letters A*, 263(4-6):341–346, dec 1999. ISSN 03759601. doi: 10.1016/S0375-9601(99)  
695 00757-4.

- 702 Carlos Oliver, Dexiong Chen, Vincent Mallet, Pericles Philippopoulos, and Karsten Borgwardt.  
703 Approximate Network Motif Mining Via Graph Learning. 2022.  
704
- 705 Mathew Penrose. *Random geometric graphs*. Oxford Studies in Probability. Oxford University  
706 Press, London, England, May 2003.
- 707 Gilles Pisier. *The Volume of Convex Bodies and Banach Space Geometry*. Cambridge University  
708 Press, October 1989. ISBN 9780511662454. doi: 10.1017/cbo9780511662454.  
709
- 710 Nataša Pržulj. Biological network comparison using graphlet degree distribution. *Bioinformatics*,  
711 23(2):e177–e183, jan 2007. ISSN 1367-4811. doi: 10.1093/bioinformatics/btl301.
- 712 Pedro Ribeiro and Fernando Silva. G-tries: a data structure for storing and finding subgraphs.  
713 *Data Mining and Knowledge Discovery*, 28(2):337–377, February 2013. ISSN 1573-756X. doi:  
714 10.1007/s10618-013-0303-4.
- 715 Pedro Ribeiro, Pedro Paredes, Miguel E. P. Silva, David Aparicio, and Fernando Silva. A survey  
716 on subgraph counting: Concepts, algorithms, and applications to network motifs and graphlets.  
717 *ACM Comput. Surv.*, 54(2), mar 2021. ISSN 0360-0300. doi: 10.1145/3433652.  
718
- 719 Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang.  
720 Self-supervised graph transformer on large-scale molecular data, 2020.
- 721 Satyaki Roy, Preetam Ghosh, Dipak Barua, and Sajal K. Das. Motifs enable communication effi-  
722 ciency and fault-tolerance in transcriptional networks. *Scientific Reports*, 10(1), June 2020. ISSN  
723 2045-2322. doi: 10.1038/s41598-020-66573-x.  
724
- 725 Shai S. Shen-Orr, Ron Milo, Shmoolik Mangan, and Uri Alon. Network motifs in the transcriptional  
726 regulation network of *Escherichia coli*. *Nature Genetics*, 31(1):64–68, 2002. ISSN 10614036. doi:  
727 10.1038/ng881.
- 728 Jinfang Sheng, Yufeng Zhang, Bin Wang, and Yaoxing Chang. MGATs: Motif-Based Graph  
729 Attention Networks. *Mathematics*, 12(2):293, jan 2024. ISSN 2227-7390. doi: 10.3390/  
730 math12020293.
- 731 Miguel E P Silva, Robert E Gaunt, Luis Ospina-Forero, Caroline Jay, and Thomas House. Compar-  
732 ing directed networks via denoising graphlet distributions. *Journal of Complex Networks*, 11(2),  
733 February 2023. ISSN 2051-1329. doi: 10.1093/comnet/cnad006.  
734
- 735 Petar Veličković, Arantxa Casanova, Pietro Liò, Guillem Cucurull, Adriana Romero, and Yoshua  
736 Bengio. Graph attention networks. *6th International Conference on Learning Representations,*  
737 *ICLR 2018 - Conference Track Proceedings*, pp. 1–12, 2018. doi: 10.1007/978-3-031-01587-8\  
738 \_7.
- 739 Petar Veličković, Rex Ying, Matilde Padovano, Raia Hadsell, and Charles Blundell. Neural Exe-  
740 cution of Graph Algorithms. *8th International Conference on Learning Representations, ICLR*  
741 *2020*, 2020.
- 742 Yanbo Wang and Muhan Zhang. Towards Better Evaluation of GNN Expressiveness with BREC  
743 Dataset. apr 2023.  
744
- 745 Shuhei Watanabe. Tree-structured parzen estimator: Understanding its algorithm components and  
746 their roles for better empirical performance, 2023.
- 747 Duncan J. Watts and Steven H. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*,  
748 393(6684):440–442, June 1998. ISSN 1476-4687. doi: 10.1038/30918.  
749
- 750 Anatol E. Wegner. Motif Conservation Laws for the Configuration Model. pp. 4–6, aug 2014.
- 751 Zhenqin Wu, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S.  
752 Pappu, Karl Leswing, and Vijay Pande. Moleculenet: A benchmark for molecular machine learn-  
753 ing, 2018.  
754
- 755 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How Powerful are Graph Neural  
Networks? pp. 1–17, oct 2019.

756 Zuoyu Yan, Junru Zhou, Liangcai Gao, Zhi Tang, and Muhan Zhang. Efficiently Counting Substructures by Subgraph GNNs without Running GNN on Subgraphs. mar 2023.  
757  
758

759 Rex Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. GNNExplainer:  
760 Generating Explanations for Graph Neural Networks. *Advances in neural information processing*  
761 *systems*, 32:9240–9251, dec 2019. ISSN 1049-5258.

762 Rex Ying, Andrew Z. Wang, and Jure Leskovec Jiaxuan You. Spminer: Frequent subgraph mining  
763 by walking in order embedding space, 2020.  
764

765 Zhaoning Yu and Hongyang Gao. Molecular representation learning via heterogeneous motif graph  
766 neural networks, 2022.

767 Bohang Zhang, Guhao Feng, Yiheng Du, Di He, and Liwei Wang. A Complete Expressiveness  
768 Hierarchy for Subgraph GNNs via Subgraph Weisfeiler-Lehman Tests. *Proceedings of Machine*  
769 *Learning Research*, 202:41019–41077, feb 2023. ISSN 26403498.

770 Bohang Zhang, Jingchu Gai, Yiheng Du, Qiwei Ye, Di He, and Liwei Wang. Beyond Weisfeiler-  
771 Lehman: A Quantitative Framework for GNN Expressiveness. pp. 1–73, jan 2024.  
772

773 Shichang Zhang, Ziniu Hu, Arjun Subramonian, and Yizhou Sun. Motif-Driven Contrastive Learn-  
774 ing of Graph Representations. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35  
775 (18):15980–15981, dec 2020. ISSN 2374-3468. doi: 10.1609/aaai.v35i18.17986.  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## 810 A DATA DETAILS

### 811 A.1 NON-DETERMINISTIC GENERATORS

812 For the Erdős-Renyi model Erdős & Rényi (1960), we mainly aim at creating graphs in the three  
813 (out of four) main topological phases a graph can be (Barabási & Pósfai, 2017). We exclusively  
814 uniformly control the number of nodes within each of the delineated phases, namely, the “critical”,  
815 “supercritical” and “connected” states. This strategic regulation facilitates substantial variability in  
816 graph size while preventing an excessive escalation of the referred metric that could possibly impede  
817 further computational processing.

818 For the Watts-Strogatz (Watts & Strogatz, 1998) and Newman Watts-Strogatz (Newman & Watts,  
819 1999), we regulated the generation based on the total number of nodes, the initial number of neigh-  
820 bours and the probability of rewiring in order to generate networks that represented key sections of  
821 the characterisation based on the clustering coefficient and path length, as given in Watts & Strogatz  
822 (1998).

823 For the extended Barabási-Albert model (Albert & Barabási, 2000), we defined as hyperparameters  
824 the total number of nodes and amount of connections a new node gains. Subsequently, we em-  
825 ploy the equations delineated in the original article, values for the probabilities associated with the  
826 formation of new links ( $p$ ) and the rewiring of existing connections ( $q$ ) are derived. These computa-  
827 tions aim to yield graphs characterised by a power-law degree distribution with an exponent ranging  
828 uniformly between 2 and 3.

829 For the cluster power-law (Holme & Kim, 2002), we vary uniformly the number of nodes and  
830 calculate the necessary probability according to the original study to obtain a clustering coefficient  
831 of 0.35, 0.45 or 0.55.

832 The duplication-divergence generator (Ispolatov et al., 2005) operates by randomly selecting a node  
833  $v$  from an initial graph and duplicating all edges connected to  $v$  with a retention probability denoted  
834 as  $\sigma$ . Two of the selected regimes exhibit self-averaging behaviour concerning the number of edges,  
835 specifically when  $0 < \sigma < e^{-1}$  and  $e^{-1} < \sigma < 1/2$ . The non-self-averaging regime is characterised  
836 by  $1/2 < \sigma < 1$ . More characteristics regarding the generated graphs, for example, the degree  
837 distribution, can be found in the original paper.

838 In the Gaussian random partition model, proposed by Brandes et al. (2003),  $k$  groups of nodes  
839 are generated with  $t$  nodes derived from a Gaussian distribution with mean  $s$  and variance  $v$ . The  
840 connectivity between nodes in a group is given by a probability  $p$ , and the connectivity inter-groups  
841 is given by  $q$ . In this generator, we parameterise the number of nodes  $|V|$ , the size of the  $k$  groups  
842 and the maximum number of allowed edges  $|E_{\max}|$ . Both the  $p$  and  $q$  probabilities are calculated to  
843 not exceed the maximum number of edges according to Equation 1.

$$844 q \leq \min\left(1, \frac{2|E_{\max}|}{|V|^2 + |V|(\kappa \cdot s^{1/2} - s(\kappa + 1))}\right) \quad (1)$$

$$845 p \leq \min(1, \kappa \cdot q) \quad (2)$$

846 We defined  $p$  as always having the possibility of going above  $q$  because we would like to have  
847 networks that can have a some community structure in order to have a more diverse set of graphs.  
848 Hence, we put the upper bound of  $p$  as being scaled over  $q$  by  $\kappa$ , which we called over-attractiveness.  
849 The values used for the  $v$  and  $\kappa$  are 10 and 5 respectively. All other parameters are uniformly  
850 sampled from a predefined range<sup>1</sup>.

851 In the case of the forest-fire model (Leskovec et al., 2007), we varied the number of nodes and  
852 the backward and forward probability between 0 and 0.4 (inclusive) to try to steer away from very  
853 aggressive Densification Power Law exponents and clique-like graphs, characteristics that, if se-  
854 vere, can hinder the subsequent steps from a computational point of view. With the values for the  
855 probabilities defined above, we expect to observe sparse networks that slowly “densify over time”,  
856 together with decreasing diameter. All the graphs are made undirected after being generated.

863 <sup>1</sup>Details for the parameters available in the supplemental material.

864 For the random geometric graph, since some properties of the graph related to its connected-  
 865 ness, such as maximum cluster size and coefficient vary with the dimension of the unit hypercube  
 866 used (Dall & Christensen, 2002; Penrose, 2003), we decided to vary the number of nodes and the  
 867 dimension of the hypercube used between 2 and 5. However, we did not efficiently explore all  
 868 possible configurations within the referred dimensions because we limited the number of edges.  
 869 Similarly to a random geometric model, we used a random geometric model in 3D with duplication  
 870 divergence (Higham et al., 2008). For this model, we followed Silva et al. (2023).

871 The last model in the non-deterministic segment is the random regular generator. In this case, the  
 872 parameters subject to uniform variation were the total number of nodes and the degree assigned to  
 873 each node, which once determined, remain constant across all nodes.

## 875 A.2 DETERMINISTIC GENERATORS

876  
 877 We complemented the graphs generated by the non-deterministic generators with smaller amounts of  
 878 graphs from deterministic generators. These generators have their network completely and without  
 879 randomness determined once their parameters are chosen.

880 The first group of deterministic generators consists of multiple types of trees. We use the binomial  
 881 tree model parameterised on its order and the balanced tree (full rary-tree) parameterised on its  
 882 height and branching factor.

883 The second group is based on modified cycles. We use the circular ladder generator, varying the  
 884 complete size of the graph and the chordal cycle (Lubotzky, 1994), also varying its complete size.

885 The third group is based on complete graphs and encompasses the barbell and lollipop graphs. The  
 886 barbell graph is made of two complete graphs of size  $k$  connected by a path of size  $m$ . The lollipop  
 887 is a barbell graph with only one complete graph and the path. In order to not complicate subsequent  
 888 steps, we carefully limited the size of the complete graphs.

889 The fourth group consists of the Dorogovtsev-Goltsev-Mendes model (Dorogovtsev et al., 2002).  
 890 This generator modulates a scale-free discrete degree distribution with exponent  $1 + \ln 3 / \ln 2$  by  
 891 using a rather simple rule: “At each time step, to every edge of the graph, a new vertex is added,  
 892 which is attached to both the end vertices of the edge.” (in Dorogovtsev et al. (2002)). We vary  
 893 the magnitude of the number of nodes and edges by changing  $n$ , resulting in  $3(3^n + 1)/2$  and  $3^{n+1}$   
 894 nodes and edges respectively.

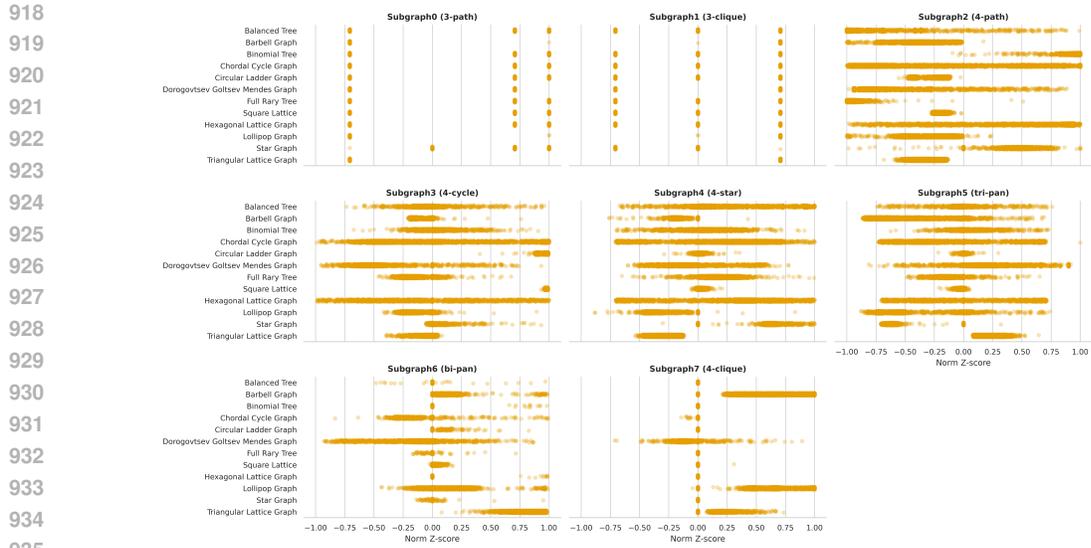
895 The fifth group consists of lattices. Namely, we use 2D hexagonal, triangular lattices and 3D square  
 896 lattices. The first 2 lattices have the option of allowing for boundary periodicity. All lattices vary in  
 897 terms of the size of each dimension.

898 Finally, the last group consists of star graphs of various sizes.

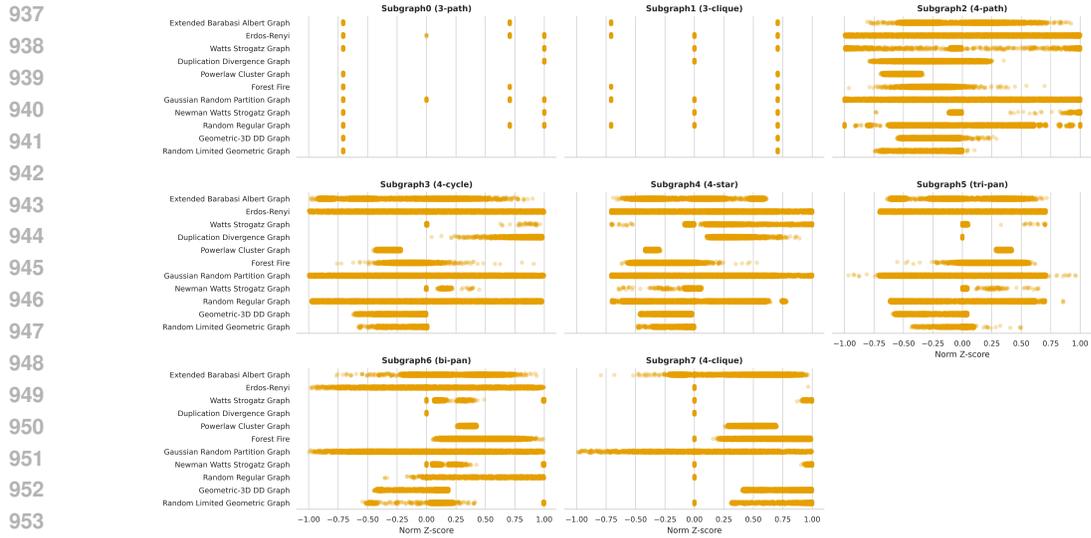
899 Since the types of graphs that the deterministic generators generate are not subject to randomness,  
 900 it is redundant to create multiple graphs for each set of parameters. However, in order to introduce  
 901 a degree of randomness to the deterministic graphs, we introduced a probability of random rewiring  
 902 of a percentage of edges after the graph is generated. The rewiring procedure for a single edge  
 903 consists of selecting an edge  $(u, v)$  from a graph  $\mathcal{G}$ , deleting it and attaching one of the ends,  $u$  or  $v$ ,  
 904 to another node  $w$ . If  $u$  is picked and  $(u, w)$  already exists, then  $\mathcal{G}$  will exit the procedure with one  
 905 less edge. Since we want some variability but still want to preserve the original deterministic graphs,  
 906 for each generator, two sets of graphs  $\mathbb{S}_1$  and  $\mathbb{S}_2$  will be generated where each set goes through all  
 907 the proposed generator parameters. After that,  $\mathbb{S}_1$  is not subject to any rewiring, and for each graph  
 908 in  $\mathbb{S}_2$ ,  $p\%$  of its edges are rewired according to the procedure described earlier. According to this  
 909 methodology, we generated four versions. The first had 2 edges swapped, the second 25% of the  
 910 edges swapped, the third 10% and the fourth 60%. We stick to version two due to being the best  
 911 performing one according to preliminary tests. This fact means that 25% seems to be a good choice  
 912 of random-rewiring so that the information encoded in the deterministic graphs is maximised.

## 914 A.3 EXACT DATASET PARTITION

915 Initially, we will conduct separate experiments using the two segments of the produced data, the  
 916 deterministic and the non-deterministic. Furthermore, the split in train-validation-test is stratified  
 917



(a) Deterministic Generators



(b) Non-Deterministic Generators

Figure 2: Distribution of the values the significance profile of each graph in  $\Omega$  takes in the synthetic dataset. Each point represents a graph.

by random sampling a percentage  $p$  of each of the generators for each segment. In the case of the deterministic segment, we use all 3200 graphs available. With  $p = 0.7$ , the training set has  $3200 \times 0.7 \times 12 = 26880$  graphs, and the validation set has  $3200 \times 0.2 \times 12 = 7680$  graphs. The remaining 10% are used for the test set. We avoided using larger datasets due to memory restrictions. As for the non-deterministic segment, in order for it to have a comparable total size to the deterministic segment, we sampled  $3490 \times 0.7 \times 11 = 26873$  graphs and the validation set  $3490 \times 0.2 \times 11 = 7678$ .

#### A.4 PATTERN INTERCONNECTIVITY

Figure 2 adopts a view of individual cases at the cost of a global view of the patterns. Each dot corresponds to an individual graph. Each point is slightly transparent to attempt to give the notion of density.

Following the discovery introduced in Section 5 regarding the symmetry of the score of size three patterns, we have a strong indication that even without the normalisation of the Z-score, the number of occurrences between connected graphs of the same size is highly related. More formally, the relation between the Z-scores of the graphs of size three can be described as follows. Let  $x$  be a random variable denoting the number of induced occurrences of triangles in any graph that follows the degree distribution  $D$ . Let  $y$  be a random variable denoting the occurrence of induced 3-paths in any graph that follows the degree distribution  $D$ . Equation 3 gives the relation between Z-scores of  $x$  and  $y$ .

$$X = \begin{cases} 0, & \text{if } y - \mu_y = 0 \\ -\frac{\sigma_x}{\sigma_y}(Y - \mu_y) + \mu_x, & \text{otherwise} \end{cases} \quad (3)$$

When standardised to a mean of 0 and a standard deviation of 1, the Z-scores of both variables, exhibit symmetry. Hence, it is possible to express their non-standardised values as linear combinations of each other. Considering the mean and standard deviation of the counts of 3-paths and triangles for  $D$ , given any graph  $\mathcal{G}$  that follows  $D$ , it is possible to get the concrete count of triangles from the count of 3-paths and vice-versa.

Even though from a practical point of view, the result from Equation 3 has little implications due to the dependence on the first raw moment and the second central moment of both the distributions of  $x$  and  $y$ , it presents a strong indication of what was postulated in Section 4 regarding the connectivity between the graphs selected for  $\Omega$ . In this case, the relation is so strong that we believe to be redundant to try to predict both scores. Moreover, following the normalisation procedure, the restrained nature of the result raises questions about the choice of modelling the problem as a regression task for the size three graphs. However, despite these observations, we stick to our initial formulation since in theory it does not undermine the ability of the model.

The result experimentally verified in the above paragraphs can be seen as a small extension of Ginoza & Mugler (2010) and Wegner (2014) to undirected patterns of size three. In particular, adapting from Wegner, Equation 4 displays the conservation law for the number of induced 3-paths.

$$\#3\text{-paths}_{ind} = \underbrace{\#3\text{-paths}_{ind}}_{\text{fully defined by degree sequence}} - \underbrace{3\#\text{triangles}}_{\text{not fully defined by degree sequence}} \quad (4)$$

Since the number of non-induced 3-paths depends only on the degree sequence ( $\sum_{i=0}^{|V|} \binom{|N(i)|}{2}$ ), it will not change under the configuration model. Hence, the number of induced 3-paths is a variable that once the degree sequence is fixed, depends only on the number of triangles. As for the number of triangles, they depend on the order the edges are added to the graph under equation 5.

$$\text{total triangles} = \sum_{t=0}^{|E|} \#\text{triangles}_t \quad (5a)$$

$$\text{total 3-path} = \sum_{t=0}^{|E|} (\#\text{3-path}_t - \#\text{triangles}_t) \quad (5b)$$

$$\#\text{triangles}_{t+1} = |\{w|w \in N(u^t) \wedge w \in N(v^t)\}| \quad (5c)$$

$$\#\text{3-path}_{t+1} = |N(u^t)| + |N(v^t)| - 2\#\text{triangles}_{t+1} \quad (5d)$$

where nodes  $u$  and  $v$  represent the nodes that were connected by an edge at iteration  $t$ . Hence, any realisation of a degree sequence through the configuration model will always have its number of induced 3-paths negatively correlated with the number of triangles.

Regarding the relation between graphs of size three and graphs of size four, by analysing Figure 3, it is possible to understand that there is a relation between the significance profiles of these graphs. This relationship is particularly pronounced concerning the 4-star, tri-pan and 4-clique, as the values

of the significance profiles assumed by these graphs are distributed across the spectrum centred at 0, contingent upon the value held by the 3-path. As for the 4-cycle and bi-fan, this relation is not as strong. For the 4-cycle, we learn that the values are mostly zero when the significance-profile for the 3-path is negative and is quite dispersed across the space otherwise. As for the bi-fan, even though hard to discern from the figure, 46.2% of the values are 0 when the significance-profile for the 3-path is positive, and 69.7% are between  $-0.1$  and  $0.1$  for the same conditions.

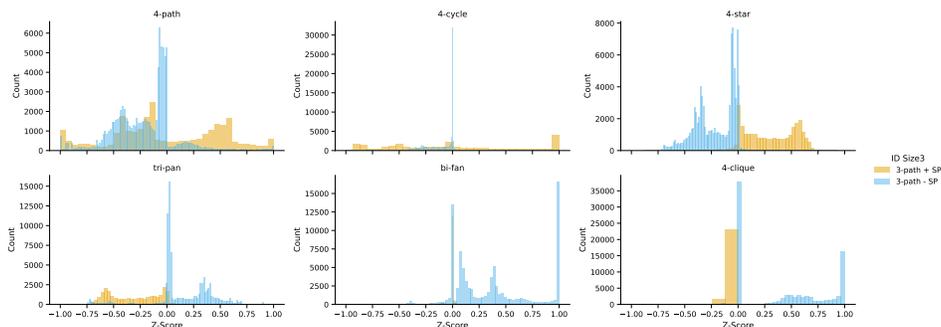


Figure 3: Distribution of the significance profiles for the graphs of size 4, given the value the 3-path took. The positive value corresponds to  $0.707106$  and the negative to  $-0.707106$ .

#### A.5 REAL-WORLD DATASET

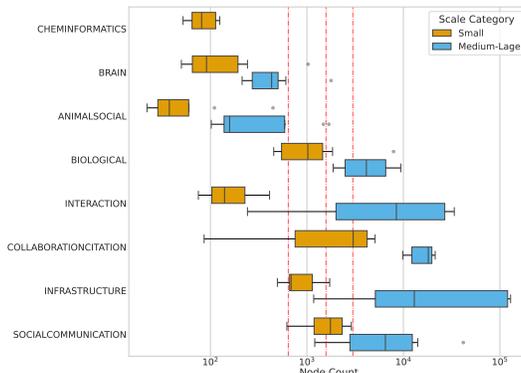
The selected type categories are described in list A.5. The numbers between the square brackets in each bullet point correspond to the number of networks each category has in each of the scale categories from the smallest to the largest scale.

- **ANIMAL SOCIAL**: [10/8] networks of the social behaviour of non-human animals incorporating a spectrum that includes ants, dolphins, lizards, sheep, and other examples.
- **BIOLOGICAL**: [10/10] networks of protein-protein interactions, a metabolic network of small organisms and a network of connections between diseases in humans by the number of shared genes.
- **BRAIN**: [9/10] networks of connectomes of diverse regions such as the cerebral cortex, interareal cortical network, and neural synaptic, among others, of multiple species such as cats, worms, mice, macaques and humans.
- **CHEMOINFORMATICS**: [10/0] networks of multiple different enzyme structures.
- **COLLABORATION CITATION** [6/8]: networks of citations of papers and collaborations between authors.
- **INFRASTRUCTURE**: [5/7] Electric grids and road networks.
- **INTERACTION**: [5/6] Networks of physical contact between humans in various contexts, together with some digital contact, for example, by e-mail or a phone call.
- **SOCIAL COMMUNICATION**: [2/10] Interaction between humans in social networks such as mutually liked Facebook pages, friendship connections and retweets.

Figure 4a and Figure 4b show a summary of the number of edges and nodes for the different type and scale categories. The red dashed lines represent the average of the minimum node (or edges) quantity, the average of the mean node (or edges) quantity and the average of the maximum node (or edges) quantity respectively, calculated for all 23 graph models used in the synthetic dataset.

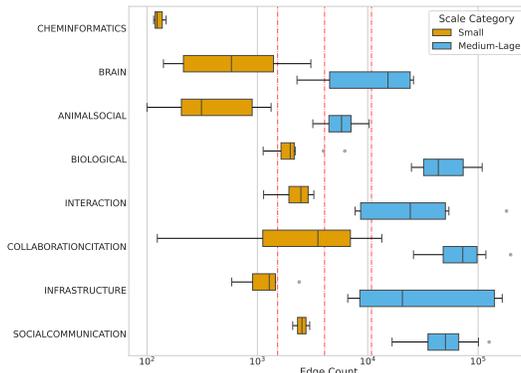
In Figure 4a and Figure 4b, a noticeable distinction is observable in the delineation of scale categories based on the type category. The distinction between scale categories is influenced by the type category of the network. For example, ANIMAL SOCIAL networks tend to be smaller than INFRASTRUCTURE networks in the medium-large category. However, this relationship varies by scale, as both network types exhibit similar sizes in the small-scale category, contrasting with their

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091



(a) Summary of how the number of nodes is distributed for the real networks.

1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105



(b) Summary of how the number of edges is distributed for the real networks.

1106  
1107

Figure 4: Summary of the distribution of the node and edge count of the real networks. All data is presented in logarithmic scale.

1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117

medium-large behaviour. In any case, in the general scenario, we expect that the small-scale category encloses networks from being smaller than an average network from the training set until networks that can be slightly larger than the mean training network. As for the medium-large category, they hold networks that can be around the average size of a training network to networks that are several orders of magnitude larger than the average size of a training network. The dashed red lines in Figure 4a and Figure 4b help validate this statement.

1118  
1119  
1120  
1121

The supplemental details have further details regarding the real-networks used. In summary, we have 56 networks in the small-scale category and 59 in the medium-large category. If a graph in the test set was not already a simple undirected static graph when it was obtained, we transformed it in a graph following said conditions.

1122  
1123  
1124  
1125

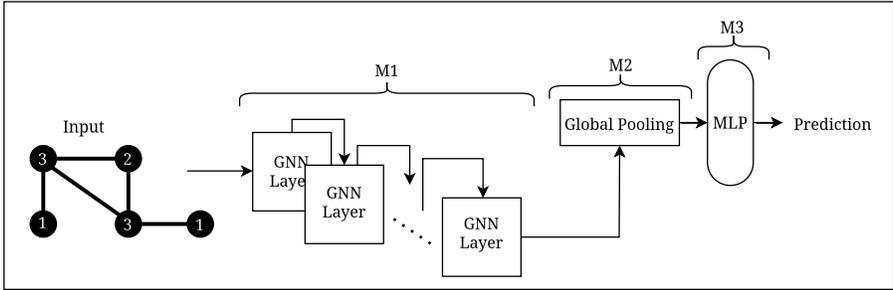
## B INITIAL BASE MODEL

1126  
1127  
1128  
1129  
1130

The model ( $\mathcal{B}$ ) consists of three modules. The first module consists of  $K$  layers of a GNN. The job of the first module is to work on the graph data and adjust the node embeddings so that the second module, a global pooling function, can summarise them into a single graph-level embedding. The third module is an MLP that takes as input the graph embedding and will adapt it to output the final prediction for the normalised Z-scores of the graphs in  $\mathcal{s}$ . Figure 5 shows a diagram of the model.

1131  
1132  
1133

All the optimisations for the hyperparameters of  $\mathcal{B}$  will be performed by Optuna (Akiba et al., 2019) with 450 rounds of suggestions of hyperparameters, orchestrated through Ray (Liaw et al., 2018; Moritz et al., 2018). Moreover, the hyperparameter sampling procedure employed the Tree-Structured Parzen Estimator (Watanabe, 2023), while the pruning strategy was executed through the

Figure 5: Illustration of the base model  $\mathcal{B}$  divided in three modules, M1, M2 and M3.Table 2: Break down of the hyperparameter space used for  $\mathcal{B}$ .

	Min	Max		Epochs	Batch Size	Learning Rate (log)
M1-GNN Depth	2	3		100	16, 32, 64, 128, 256	[0.00001, 0.001]
M1-Hidden Dimension	6	16				
M1-GNN Dropout	0.0	0.9				
M1-Jumping Knowledge	max, cat, lstm					
M2-Global Pool		add				
M3-MLP Depth	2	6				
M3-Hidden Dimension	6	16				
M3-MLP Dropout	0.2	0.9				

application of the median rule (Golovin et al., 2017). Table 2 presents the hyperparameter space used for model  $\mathcal{B}$ .

The asymmetry in the hyperparameter space presented in Table 2 stems from our choice of preemptively test a slightly larger hyperparameter space and identify some values that resulted in very bad results. From this early testing phase, we also narrowed down M3 from a global add, mean or max function to just the global add function. This result aligns with some limitations that are known for the mean and max pooling functions (Xu et al., 2019). Furthermore, the fixed values of 100 epochs can be shortened not only by the pruner but also by an early-stopping module with a grace period of 25 epochs, synced with the median pruner, and patience of other 25 epochs of not seeing an improvement for the global minimum loss. Moreover, since we believe our problem does not need very long range dependencies since the structures in  $\Omega$  can be fully defined by a hop size of 2, in order to try to limit the problem of over-smoothing we limited the maximum number of GNN layers to 3 based on the findings that most networks have a small diameter (Albert et al., 1999; Barabási et al., 2000; Watts & Strogatz, 1998). By limiting the GNN layers, we also hope to reduce over-squashing.

For the initial tests, we performed several experiments on the synthetic dataset described in Section 5 for different M1 modules. Each experimental iteration comprised 450 trials, each uniquely characterised by a distinct combination of hyperparameters suggested by Optuna.

The four different GNNs used were GAT (Veličković et al., 2018), SAGE (max-pooling) (Hamilton et al., 2017), GCN (Kipf & Welling, 2017) and GIN (Xu et al., 2019) and are all 1-WL limited. Most of the training was done using a single NVIDIA RTX 3090 and later a NVIDIA RTX A6000.

## C METHODOLOGIES

**Persistent Patterns.** To evaluate the impact of a value  $t$  in the number of persistent patterns, we employ agglomerative clustering with complete linkage over the significance profiles with  $t$  as a stopping point for agglomeration. Using the complete linkage over  $d_\infty$  ensures that all significance profiles in a cluster remain within  $t$  of each other, effectively counting the number of persistent patterns through the number of clusters and their persistency through the size of each cluster. To discover  $t$ , we test different values and iteratively evaluate their impact on the quantity of patterns induced. The selected value is the lowest one that produces the largest drop in the number of per-

Table 3: Summary of the number of persistent patterns and their persistency. The number of clusters gives the number of persistent patterns and the statistics about their size pertain to  $\rho_s$ .

Segment	Model	Q1	Q2	Q3	Min	Mean	Max	STD	Number Clusters	Number of SPs	Threshold
ND - Test	True	1.00	1.00	1.00	1.00	1.65	56.00	3.15	2325	3839	0.01
	GIN	1.00	1.00	2.00	1.00	19.39	710.00	73.92	198		
	SAGE	1.00	1.00	2.00	1.00	11.53	751.00	63.79	333		
D - Test	True	1.00	1.00	1.00	1.00	1.90	38.00	3.25	1935	3678	0.01
	GIN	1.00	1.00	2.00	1.00	4.73	520.00	25.88	777		
Medium-Large	True	1.00	1.00	1.00	1.00	1.38	5.00	0.94	42	58	0.08
	GIN-ND	1.25	2.00	5.25	1.00	4.14	20.00	5.07	14		
	SAGE-ND	1.00	1.00	3.00	1.00	2.15	6.00	1.61	27		
	GIN-D	1.00	1.00	2.00	1.00	2.90	27.00	5.86	20		
Small	True	1.00	1.00	2.00	1.00	1.51	7.00	1.19	37	56	0.21
	GIN-ND	1.00	1.00	3.00	1.00	2.95	12.00	3.41	19		
	SAGE-ND	1.00	2.00	3.50	1.00	2.95	10.00	2.57	19		
	GIN-D	1.00	1.00	3.00	1.00	2.67	14.00	3.43	21		

sistent patterns (similar to the elbow method for clustering). Table 3 conveys summary statistics regarding the discovered persistency of patterns.

**Red Line.** Equation 6 delineates the derivation of the expected loss for the baseline depicted as a red line. The first step in said equation comes from expanding the square and the second step from the definition of variance for a standard uniform distribution. Conceptually, the red line embodies the unequivocal baseline for any model endeavouring to address the problem of predicting normalised Z-scores as delineated in Section 4. Any model falling short of the benchmark set by the red line can be confidently deemed inadequate.

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}((y_i - \hat{Y}_i)^2) \tag{6a}$$

$$\equiv \frac{1}{n} \sum_{i=1}^n (y_i^2 + Var(\hat{Y}_i)) \tag{6b}$$

$$\equiv \frac{1}{n} \sum_{i=1}^n (y_i^2 + 1/3) \tag{6c}$$

$$\equiv 1/3 + 2/n \tag{6d}$$

**Blue Line.** The derivation is similar to the one employed for the red line, hence omitted. Compared to the red line, the blue line defines an improved standard that any model that tries to predict the normalised Z-score as postulated in Section 4 must also clear. To delineate the value of the blue line, we formulated a model enabling the stochastic prediction of values that adhere to the constraint that the sum of the squared values of the scores of the graphs of each group of  $\Omega$  must be equal to 1. We chose to articulate the model using a set of independent standard normal variables  $\mathbb{S} = \{x_1, \dots, x_n\}$ , which are then normalised by the Euclidean norm of the respective  $\mathbb{S}$ . This particular formulation seems to manifest a good distribution across the space formed by  $\mathbb{S}$  in terms of uniformity, especially when considering the increasingly improbable nature, attributed to both the concentration of measure phenomenon and Dvoretzky’s theorem (Pisier, 1989; Giannopoulos & Milman, 2000), of truly attaining a random uniform distribution of points across the entire volume of a compact, symmetric, convex subset within an  $n$ -dimensional Banach space, like the  $n$ -Euclidean space,  $S$ , delineated by the process used to generate random guesses. On another note, precisely due to this observation, we conjecture that the problem of interest is ill-posed in very high dimensional spaces, meaning it becomes increasingly hard to have a meaningful significance-profile and thus predict them as the groups of  $\Omega$  increase in cardinality.

**Light-Pink Bar.** To ascertain the referred mean error that the light pink bar depicts, we randomised the weights of each model a total of 100 times and predicted the significance-profiles for the real-world dataset for each randomisation. In conjunction with the mean error estimation, we provide an error bar indicative of one standard deviation.

## D EXPERIMENT RESULTS

Figure 6 and 7 shows the summary of the results from the 450 rounds of hyperparameter optimisation for each model used in M1. The solid line represents the mean score, and the semi-transparent bound around each line represents the standard error. The displayed metrics are the MSE for the train and validation data, the median absolute error,  $med(\{med_i(|y_i - \hat{y}_i|, \forall i \in |y|)\})$ , the maximum absolute error calculated for a full prediction of a significance profile, and the mean value for the worst-performing prediction of a graph from  $\Omega$ . The maximum error is given by  $max(\{\sum_{j \in |\Omega|} |Y_{[i,j]} - \hat{Y}_{[i,j]}|, \forall i \in |D_{valid}|\})$  where  $Y$  is a 2-d matrix where the first dimension gives the number of examples in the validation dataset and the second dimension the length of  $s$ . As for the mean value of the worst-performing predictions, it is given by  $mean(max\{\sum_{i \in |D_{valid}|} |(Y_{[i,j]} - \hat{Y}_{[i,j]})|, \forall j \in |\Omega|\})$ .

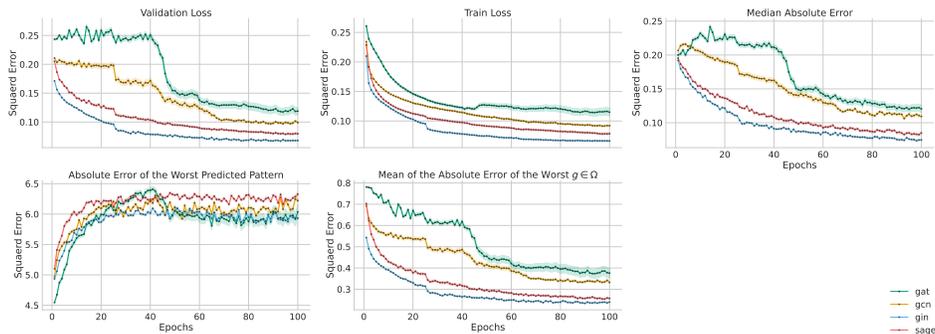


Figure 6: Learning curves for the various backends used for M1 when trained with the deterministic segment of graph generators.

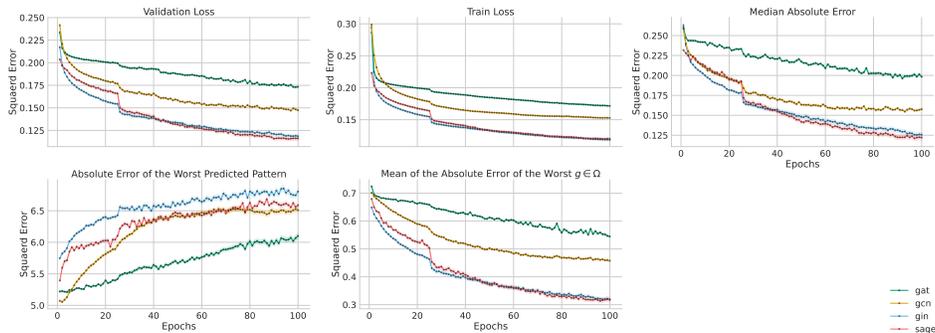


Figure 7: Learning curves for the various backends used for M1 when trained with the non-deterministic segment of graph generators.

The learning curves for the deterministic segment (Figure 6) show that all models improve significantly within the first 50 epochs, especially in all metrics except for the maximum absolute error. GIN outperforms all other models by a wide margin, prompting its selection for further analysis. For the non-deterministic segment (Figure 7), the performance of GraphSAGE and GIN is very close, with GraphSAGE holding a slight numerical edge. Since both models perform comparably, both will be retained for further evaluation.

### D.1 VALIDATION OF THE ASSUMPTIONS MADE

The two main assumptions by us proposed regards to using multi-target regression and directly predicting the significance-profiles of the chosen graphs.

**Single-target vs. Multi-target.** To allow this comparison, we trained eight models, each for one of the graphs in  $\Omega$ . The data used was the non-deterministic synthetic dataset. For the task, we utilised

Table 4: Percentiles in the validation set of the squared error and their percent decrease/increase comparing the multitarget to the single target model. Order: 3-path (0), triangle (1), 4-path (2), 4-cycle (3), 4-star (4), tailed-triangle (5), 4-chord-cycle (6), 4-clique (7).

Graph	7		6		5		4		3		2		1		0	
Type	single	multi	single	multi	single	multi	single	multi	single	multi	single	multi	single	multi	single	multi
100%	1.835	2.868	1.830	1.126	1.126	0.816	1.498	0.913	1.056	1.268	2.275	2.388	1.814	2.193	2.940	2.621
	+56.294%		-38.470%		-27.531%		-39.052%		+20.076%		+4.967%		+20.893%		-10.850%	
95%	0.294	0.250	0.376	0.334	0.304	0.320	0.350	0.348	0.547	0.502	0.517	0.489	1.294	0.741	1.463	0.808
	-14.966%		-11.170%		+5.263%		-0.571%		-8.227%		-5.416%		-42.736%		-44.771%	
75%	0.080	0.091	0.085	0.050	0.052	0.041	0.083	0.053	0.067	0.038	0.082	0.100	0.210	0.042	0.357	0.048
	+13.750%		-41.176%		-21.154%		-36.145%		-43.284%		+21.851%		-80.000%		-86.555%	
50%	0.007	0.009	0.031	0.013	0.005	0.006	0.008	0.007	0.005	0.003	0.007	0.012	0.051	0.007	0.042	0.007
	+28.571%		-58.065%		+20.000%		-12.500%		-40.000%		+71.429%		-86.275%		-83.333%	
25%	0.000	0.001	0.002	0.001	0.001	0.000	0.001	0.001	0.001	0.000	0.000	0.002	0.005	0.000	0.012	0.003
	+100.000%		-50.000%		-100.000%		0.00%		-100.000%		+100.000%		-100.000%		-75.000%	

Table 5: Percentiles in the test set of the absolute difference between the true and predicted significance-profile by direct estimation (SP) and their percent decrease/increase comparing to the predictions using a multi-target model with graph frequencies as output. Order: 3-path (0), triangle (1), 4-path (2), 4-cycle (3), 4-star (4), tailed-triangle (5), 4-chord-cycle (6), 4-clique (7).

Graph	7		6		5		4		3		2		1		0	
Type	Count	SP	Count	SP	Count	SP	Count	SP	Count	SP	Count	SP	Count	SP	Count	SP
75%	0.806	0.222	0.248	0.323	0.150	0.199	0.120	0.227	0.073	0.173	0.541	0.298	0.362	0.172	0.412	0.278
	-72.429%		+30.320%		+33.280%		+89.686%		+137.311%		-51.763%		-68.122%		-32.566%	
50%	0.161	0.116	0.046	0.109	0.065	0.056	0.065	0.083	0.042	0.083	0.202	0.126	0.476	0.072	0.335	0.079
	-28.068%		+140.276%		-14.035%		+28.404%		+99.939%		-37.503%		-84.891%		-76.506%	
25%	0.081	0.044	0.025	0.042	0.016	0.036	0.019	0.031	0.018	0.020	0.083	0.041	0.396	0.027	0.231	0.021
	-46.267%		+65.427%		+127.724%		+61.485%		+9.853%		-50.814%		-93.296%		-90.737%	

the GIN variant of MPNNs, as it demonstrated both theoretical and practical superiority among available options. The models were trained without any prior assumptions; they were initialised with the same hyperparameter space as all other models, allowing the optimiser to explore the entire parameter space. Consequently, the models are designed to specialise in their respective graph. Table 4 shows the percentiles for the squared difference between the true and predicted (using GIN model from Figure 7) significance-profile in the validation dataset for each of the eight graphs of  $\Omega$ . Each percentile has a percent comparison between the results of the multiple models.

Following the results from Table 4, apart from graph 7 (four-node clique) and 2 (four-path), generally, all graphs show an improvement in the predicted score when using multi-target regression. This result confirms our expectation delineated in Section 4 that graphs that predicting together graphs that share common traits can improve the outcome. Regarding the increased error observed in graphs 7 and 2, we hypothesise that this may be due to their limited benefit from shared information. Other graphs do not possess enough encoded information to be leverage by the shared knowledge to outperform a specialised model. Regardless, having into account the magnitude of the increase/decrease of the the errors, we believe that multi-target is superior for motif estimation. One other added benefit is the need to train only a single model that maintains competitive training times comparing to single-target regression.

**Estimating Frequency vs. Estimating Significance-Profiles.** To allow this comparison, we trained a model to estimate the frequency of the graphs of  $\Omega$  directly. The data used was the non-deterministic synthetic dataset. For the task, we utilised the GIN variant of MPNNs, as it demonstrated both theoretical and practical superiority among available options. The model was trained without any prior assumptions; it was initialised with the same hyperparameter space as all other models, allowing the optimiser to explore the entire parameter space. Consequently, the model is designed to specialise in frequency estimation. Table 5 shows the percentiles for the absolute difference between the true and predicted significance-profile, by either direct estimation (SP) or frequency estimation (Count) in the test dataset for each of the eight graphs of  $\Omega$ . Each percentile has a percentage comparison between the outcomes of the two models.

To circumvent the need to generate 500 random networks according to NULL for each network in the test set and subsequently estimate the frequency of each subgraph within them, we opted for an estimation of the predicted Z-Scores. Let  $x$  represent a random variable corresponding to the frequency of subgraphs in a given degree-distribution as estimated by the trained model. We decompose  $x$  and its associated Z-score calculation into two variables,  $y$  and  $z$ , as shown in Equation 7.

The variable  $y$  denotes the actual frequency of a subgraph, while  $z$  captures the error introduced by the model’s approximation.

$$\hat{Z} = \frac{x - \mathbb{E}[x]}{\text{Var}[x]^{1/2}} \quad (7a)$$

$$\hat{Z} = \frac{(y + z) - \mathbb{E}[y + z]}{(\text{Var}(y)^2 + \text{Var}(z)^2)^{1/2}} \quad (7b)$$

Assuming that the difference between a value  $z \sim \mu_z$  and  $\mu_z$  is proportional to  $\sigma_z$ , minding signal indetermination, Equation 7 originates Equation 8.

$$\hat{Z} = \frac{(y - \mathbb{E}[y]) + (z - \mathbb{E}[z])}{(\text{Var}(y)^2 + \text{Var}(z)^2)^{1/2}} \quad (8a)$$

$$\hat{Z} = \frac{(y - \mathbb{E}[y]) \pm \sigma_z}{(\text{Var}(y)^2 + \text{Var}(z)^2)^{1/2}} \quad (8b)$$

All values in Equation 8 are known since they were either acquired during the training of the model (values regarding  $z$ ) or were collected during the dataset construction (values regarding  $y$ ).

By approximating the Z-scores predicted by models that perform frequency estimation, and subsequently normalising them as described in Section 4, we are able to compare these approximations with the outputs of models that directly predict the normalised Z-scores (significance profiles). We employ the percentiles of the absolute differences between the true significance profiles and those directly estimated by the models, as well as between the true significance profiles and those estimated via Equation 8. Table 5 presents this comparison. The values under “Count” correspond to the minimum difference (hence worst case comparison) resulting of all valid signal combinations according to Equation 8b.

Following the results from Table 5, we conclude that, generally, predicting significance-profiles directly improves the scores for the 3-path (graph 0), triangle (graph 1), 4-path (graph 2) and the 4-clique (graph 7). All others present a general score that is worse when directly predicting the significance-profile. The graphs that experienced overall improvement are those characterised by significant variation in their number of occurrences, depending on the type and size of the network. This result likely arises from the model’s enhanced ability to handle the extreme differences in the frequency of subgraph occurrences in networks with substantial size and topological disparities. Despite, the deteriorating the results for graphs 3 through 6, apart from 75% for graph 3 and 50% for graph 6, the increase does not magnify the order of magnitude of the errors significantly. Furthermore, the decrease for the other graphs is larger and more significant. Hence, we believe predicting significance-profiles directly to be an overall improvement to predicting graph frequencies, in the context of motif estimation for the chosen null model. Furthermore, predicting significance-profiles directly is much cheaper computationally from a motif calculation point of view.

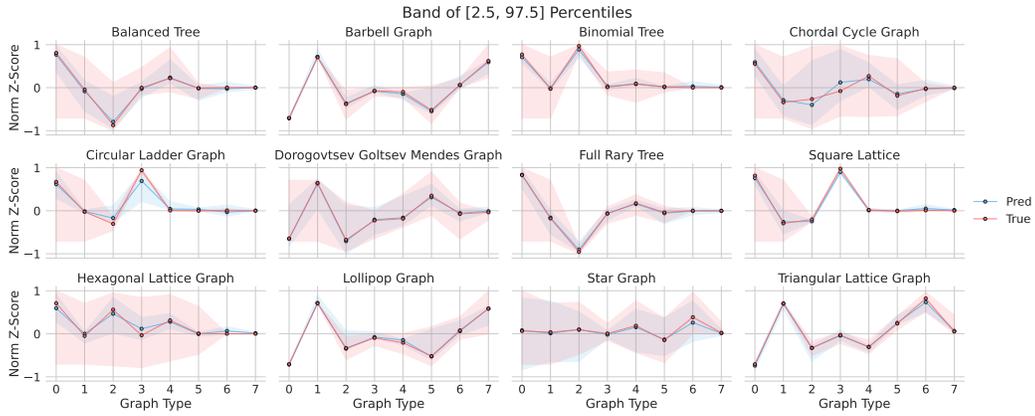
Joining the result from Table 4 and 5, in the context of motif estimation, we believe that using multi-target regression to predict significance-profiles directly outperforms using models to perform subgraph frequency estimation using either single or multi-target regression.

## D.2 PREDICTIONS

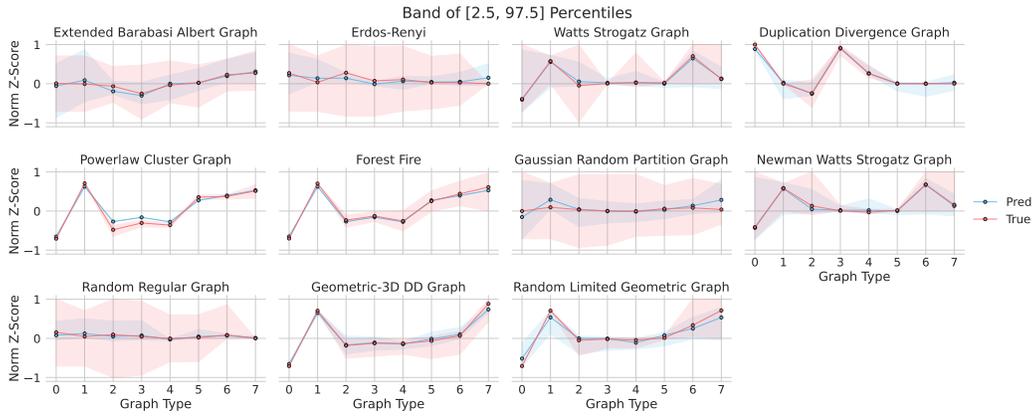
Figures 8 display a summary of the predictions for each generator made by each selected model.

Regarding the real-world dataset, since the number of graphs per category is much smaller than in the synthetic dataset, basing an analysis solely on the mean profile and percentile band can be deceiving. Regardless, we make available such figure available (Figure 12 and 13). Conversely, that scarcity allows for a more comprehensive individual analysis. Still, the volume of graphs is too high to present all images. Figures 9a through 11d show some examples. To see all predictions follow the README in the supplemental material.

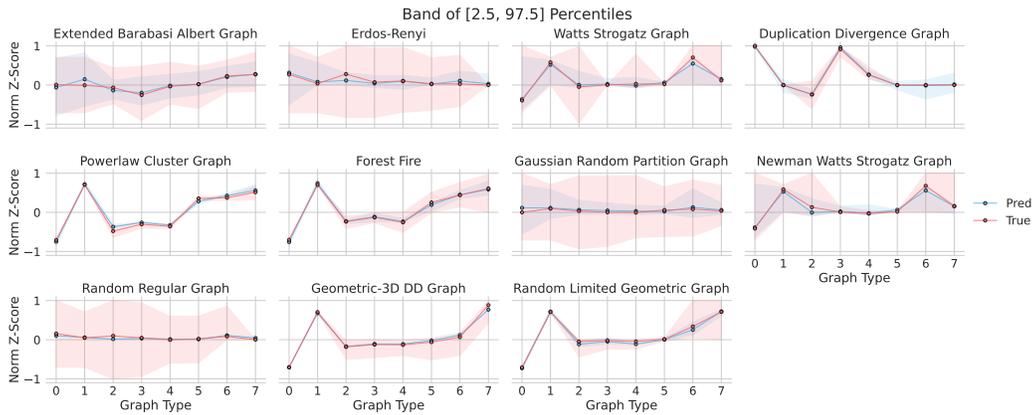
1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457



(a) GIN trained on the deterministic segment.



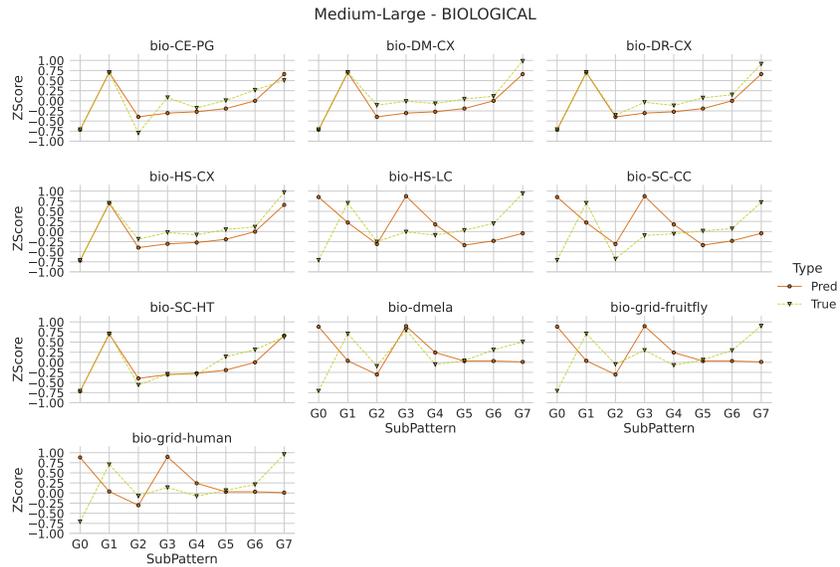
(b) GIN trained on the non-deterministic segment.



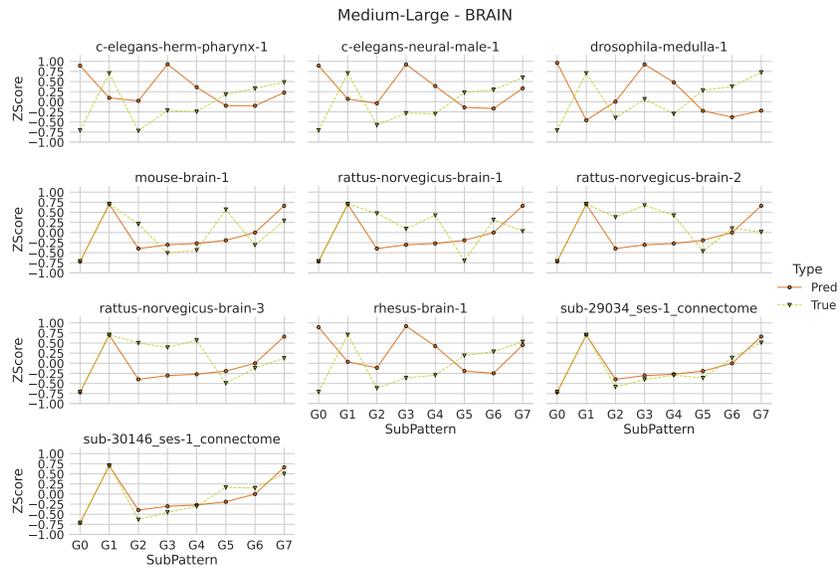
(c) SAGE trained on the non-deterministic segment.

Figure 8: Predictions for each model in each of their corresponding synthetic test datasets.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511



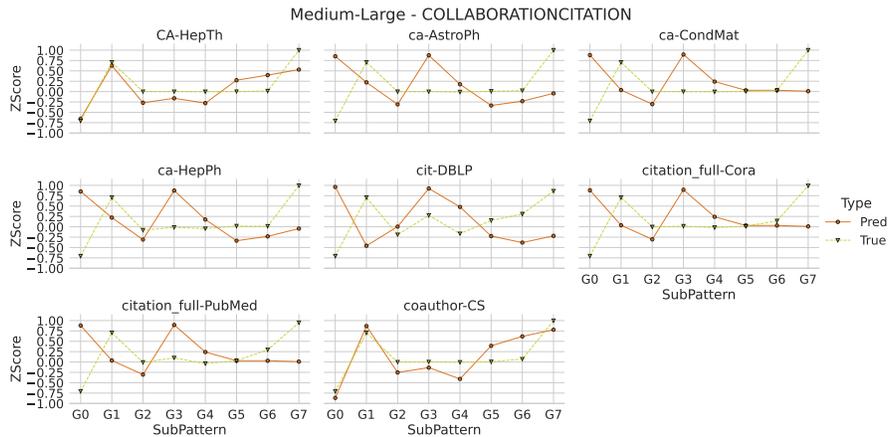
(a)



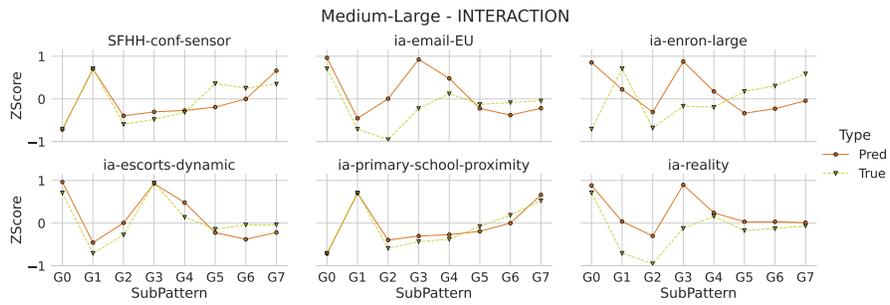
(b)

Figure 9: Continued on next page.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565



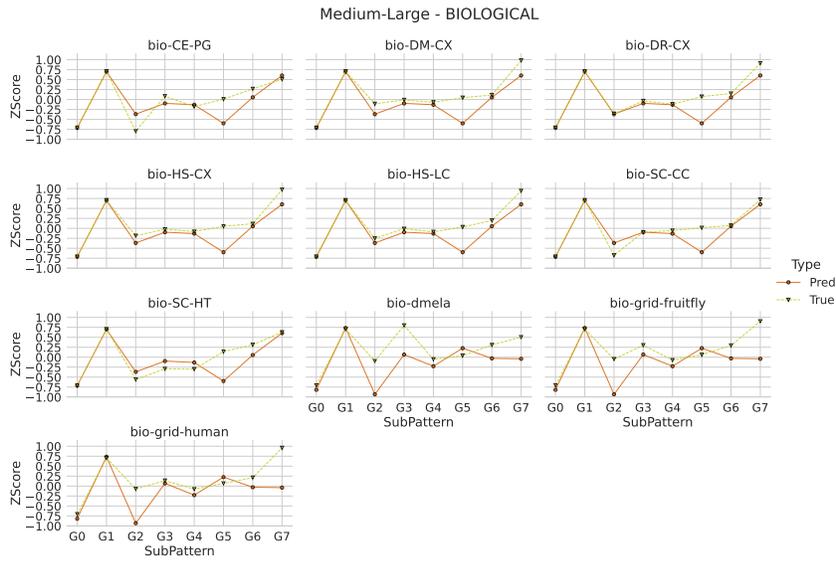
(c)



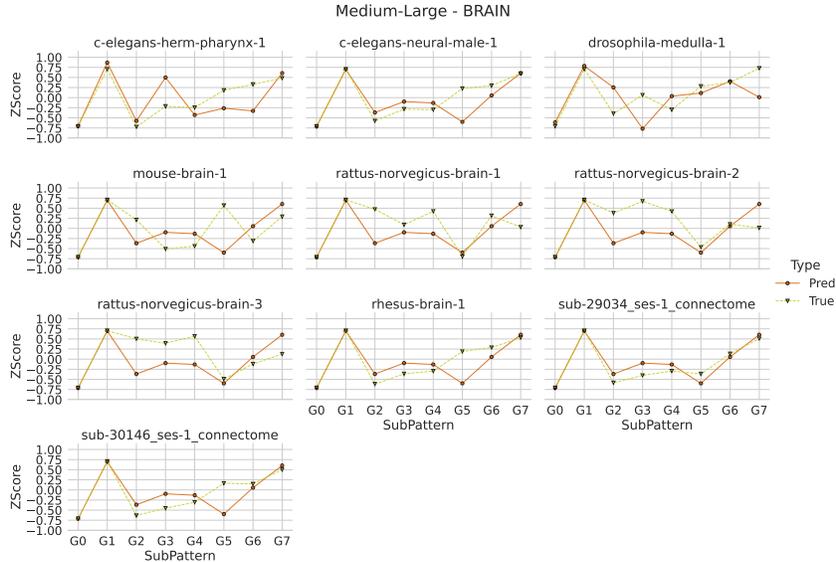
(d)

Figure 9: Predictions by GIN trained on the non-deterministic segment. Orange lines with circles are predictions and dark-yellow with triangles true values.

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619



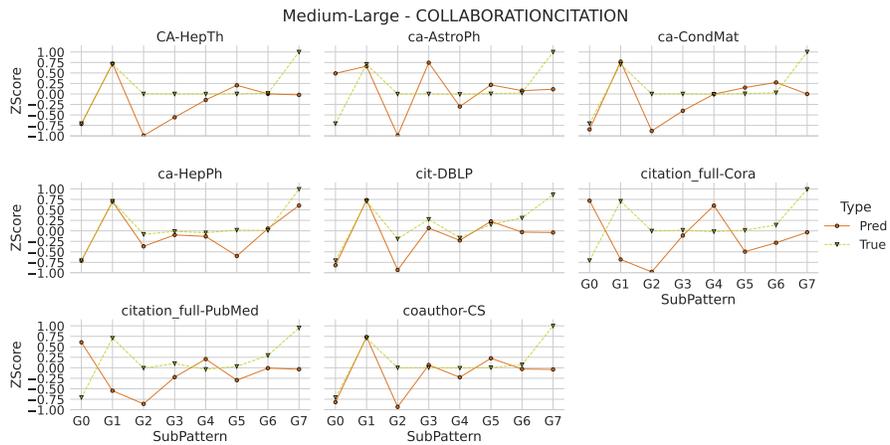
(a)



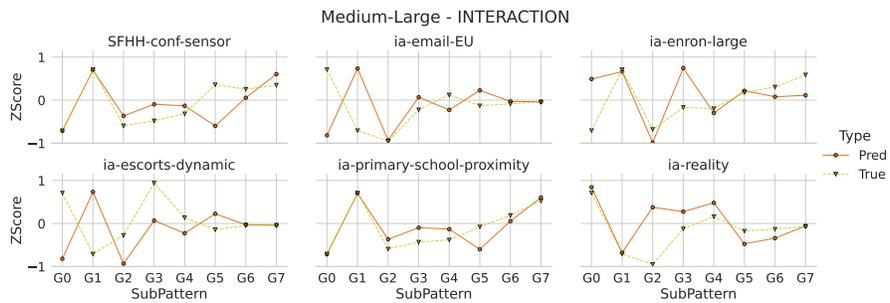
(b)

Figure 10: Continued on next page.

1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673



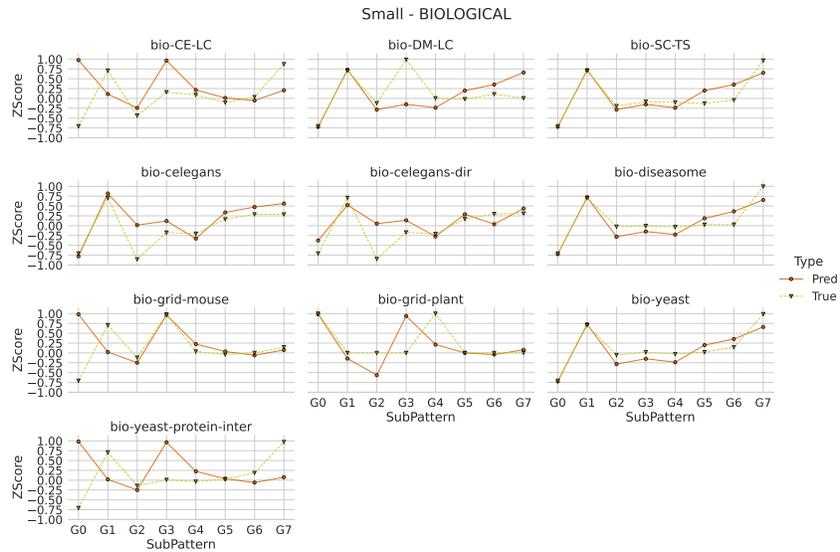
(c)



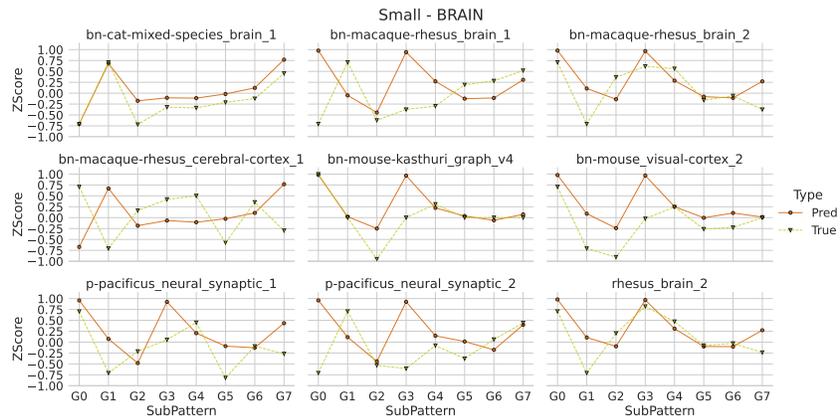
(d)

Figure 10: Predictions by GIN trained on the deterministic segment. Orange lines with circles are predictions and dark-yellow with triangles true values.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727



(a)



(b)

Figure 11: Continued on next page.

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

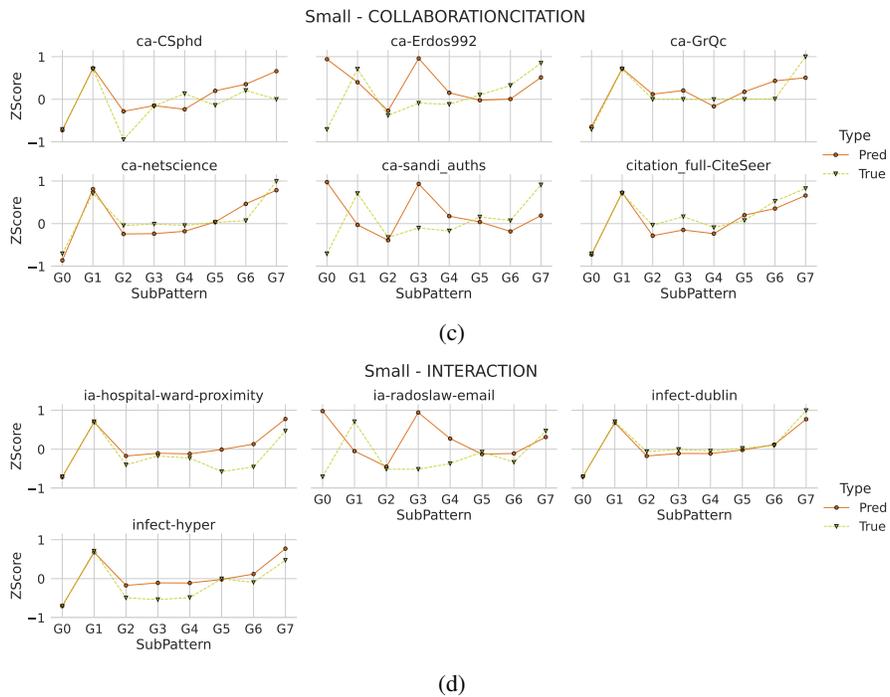
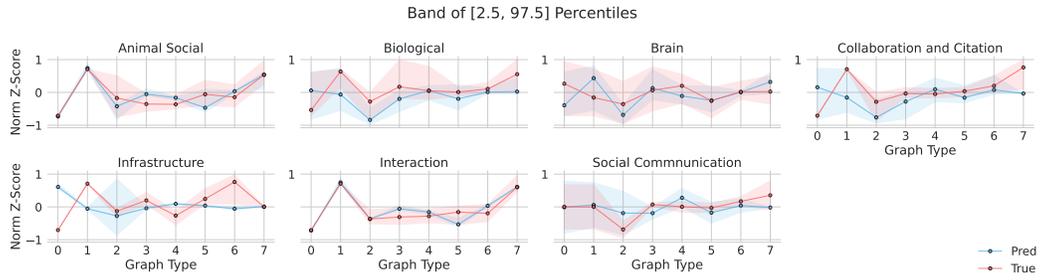
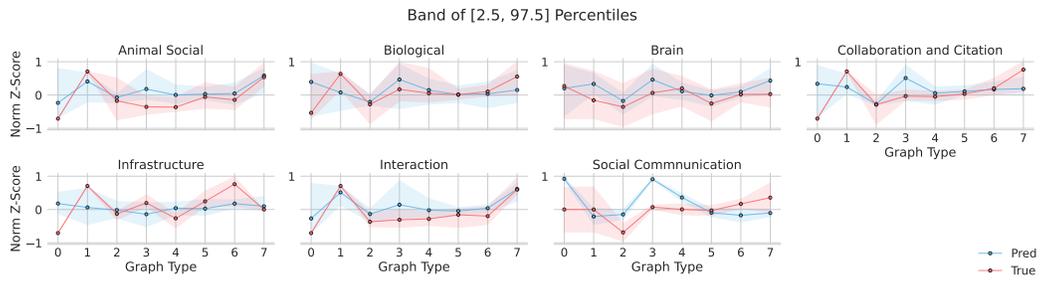


Figure 11: Predictions by GIN trained on the deterministic segment. Orange lines with circles are predictions and dark-yellow with triangles true values.

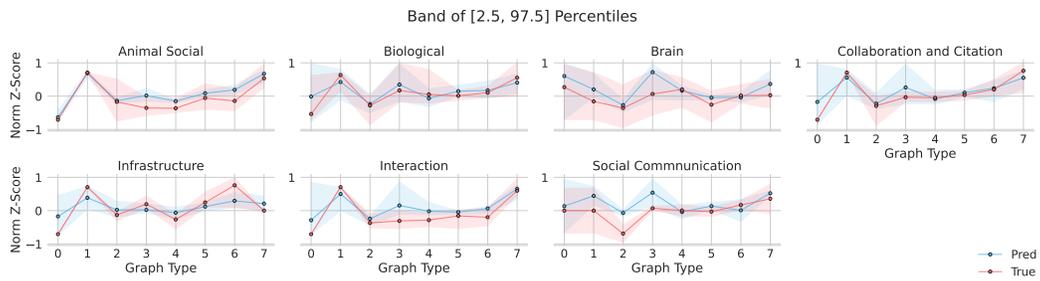
1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835



(a) GIN trained on the deterministic segment.



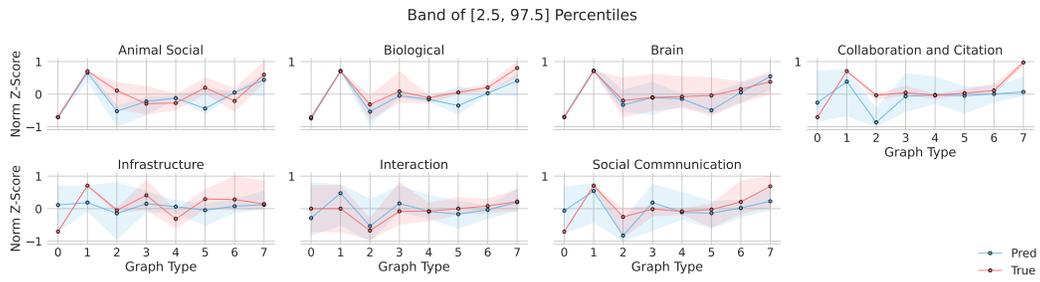
(b) GIN trained on the non-deterministic segment.



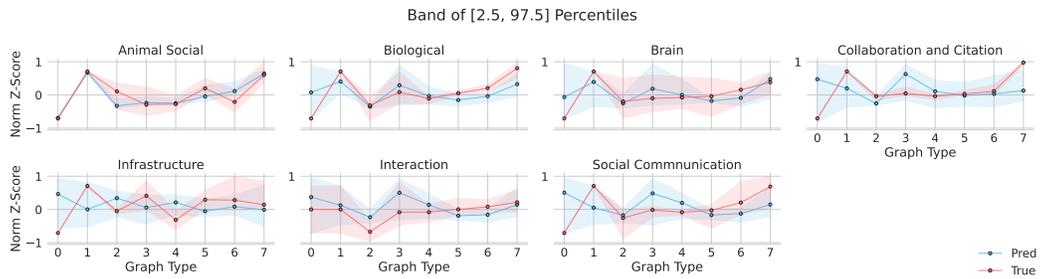
(c) SAGE trained on the non-deterministic segment.

Figure 12: Predictions for each model in for the small real-world dataset.

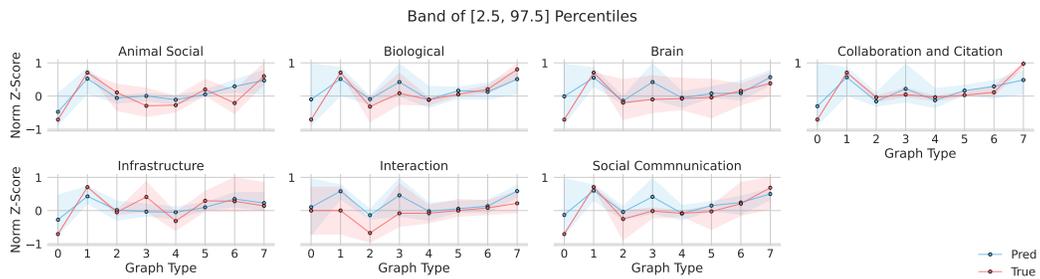
1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889



(a) GIN trained on the deterministic segment.



(b) GIN trained on the non-deterministic segment.



(c) SAGE trained on the non-deterministic segment.

Figure 13: Predictions for each model in for the medium-large real-world dataset.