



Beyond Local Descriptors: Exposing the Expressive and Computational Limits of Graph Evaluation

H. Mousavi Fard^{*1} and N. Bagherpour^{†1}

¹School of Engineering Science, College of Engineering, University of Tehran, Tehran, Iran

ABSTRACT

The evaluation of graph generative models currently relies on a fragmented ecosystem of distance metrics that fundamentally misalign with advanced likelihood-based training objectives. This paper presents a systematic critique and adversarial benchmark of the prevailing evaluation paradigms. We demonstrate that standard local statistical descriptors, such as Degree MMD, suffer from severe structural colorblindness, failing to penalize macro-structural collapse. Furthermore, metrics reliant on learned representations via Graph Neural Networks (GNNs) are theoretically capped by the 1-Weisfeiler-Lehman isomorphism test, rendering them incapable of differentiating distinct global topologies. While exact spectral methods utilizing Laplacian eigenvalues resolve these expressive limitations, our scaling analysis proves they introduce an intractable $\mathcal{O}(N^3)$ computational bottleneck. Finally, we expose the extreme statistical variance of modern classifier-based discrepancy metrics in small-sample regimes. By isolating theoretical and algorithmic failure modes, this study establishes the critical necessity for a paradigm shift toward continuous, scalable, and geometry-aware evaluation frameworks in constrained graph generation.

Keyword: Graph Generative Models, Evaluation Metrics, Spectral Graph Theory, Maximum Mean Discrepancy.

AMS subject Classification: 68R10.

^{*}Email: hesam.m.fard@ut.ac.ir

[†]Corresponding author: N. Bagherpour. Email: negin.bagherpour@ut.ac.ir

ARTICLE INFO

Article history:

Research paper

Received 02, June 2026

Accepted 13, June 2026

Available online 19, June 2026

1 Introduction

Deep generative models have rapidly advanced the capability to learn and synthesize complex structured data. In the domain of network science, architectures such as discrete diffusion models and continuous normalizing flows are optimized to reconstruct exact topological distributions. Despite these algorithmic strides in generation, the evaluation of the synthesized graphs remains a profound methodological vulnerability. The current standard for assessing graph quality predominantly relies on two-sample tests applied to heuristically selected features. This discrepancy between advanced likelihood-based training objectives and rudimentary evaluation metrics creates a persistent illusion of progress in the field.

The most prevalent evaluation framework utilizes Maximum Mean Discrepancy (MMD) to compare distributions of local topological statistics. Researchers routinely measure disparities in node degrees, clustering coefficients, and small subgraph motif counts. While computationally efficient, these local metrics are inherently susceptible to structural colorblindness [1]. A generative model can easily minimize the MMD on degree distributions while simultaneously generating graphs that are fundamentally fragmented or physically invalid. Local statistical aggregations inherently discard the long-range connectivity and global geometric constraints that define valid graph manifolds.

Attempts to address this blindness by leveraging learned representations introduce equally severe theoretical constraints. Metrics utilizing pre-trained Graph Neural Networks (GNNs) map discrete topologies into continuous latent spaces to compute distances. However, the message-passing paradigm central to standard GNNs operates as a low-pass filter, resulting in the oversmoothing of high-frequency structural anomalies. Furthermore, the expressive power of these networks is strictly bounded by the 1-Weisfeiler-Lehman (1-WL) isomorphism test [2]. Consequently, GNN-based evaluators map topologically distinct regular graphs to identical continuous vectors, assigning high validity scores to structurally corrupted outputs.

Transitioning to exact spectral evaluation resolves the topological blindness but introduces an insurmountable algorithmic bottleneck. Formulating distances based on the exact eigenvalues of the normalized Laplacian matrix captures global bottlenecks and connectivity perfectly. Yet, direct matrix diagonalization requires $\mathcal{O}(N^3)$ time complexity [3]. This cubic scaling renders direct spectral evaluation completely intractable for large-scale or dense networks. Alternative discriminative approaches, which train auxiliary classifiers to approximate Jensen-Shannon distances, mitigate the computational cost but exhibit extreme variance and overfitting tendencies in small-sample regimes.

This paper provides a systematic benchmark and critical review of the prevailing evaluation metrics for graph generative models. We aim to objectively quantify the failure modes of local descriptors, learned embeddings, and exact spectral methods. By subjecting these metrics to controlled, structure-preserving perturbations, we isolate their theoretical limits and computational boundaries. Establishing these specific failure states is a necessary precursor for developing scalable, geometry-aware evaluation frameworks.

2 Taxonomy of Existing Evaluation Metrics

The landscape of graph evaluation metrics is characterized by a fundamental tradeoff between computational tractability and structural expressivity. Existing methodologies attempt to quantify the distance between real and generated graph distributions by analyzing specific topological features. We categorize the prevailing metrics into four distinct paradigms, detailing their mathematical formulations and inherent theoretical limitations.

2.1 Local Statistical Descriptors

The most ubiquitous approach relies on comparing one-dimensional aggregations of local features. This paradigm typically utilizes the Maximum Mean Discrepancy (MMD) to evaluate the divergence between probability distributions in a Reproducing Kernel Hilbert Space. Researchers routinely compute MMD across node degree distributions, clustering coefficients, and orbit counts [4].

While computationally lightweight, these descriptors suffer from severe structural colorblindness. A generative model can reconstruct exact local degree moments while failing entirely to capture the global topological manifold. For instance, a network fragmented into multiple disconnected components can yield an identical degree distribution to a fully connected graph. Local statistics are fundamentally incapable of validating long-range connectivity constraints.

2.2 Graph Kernel Approximations

To capture extended neighborhood structures, graph kernels measure similarities through iterative subtree aggregations. The Weisfeiler-Lehman (WL) kernel is widely adopted due to its linear time complexity relative to the number of edges. It iteratively hashes neighborhood labels to construct structural signatures for each node [5].

Despite their speed, WL-based metrics are strictly bounded by their algorithmic nature. The continuous feature aggregation mirrors the 1-WL isomorphism test. Consequently, the kernel fails to distinguish between highly symmetric but topologically distinct graphs, such as strongly regular graphs with identical local symmetries but disparate global structures. It is completely blind to large cycles and global bottlenecks.

2.3 Learned Representations and GNNs

Deep representation metrics map discrete graphs into continuous latent vectors. Metrics such as the Fréchet Graph Distance (FGD) feed the synthesized networks through a pre-trained Graph Neural Network (GNN) and compute the Fréchet distance on the resulting embeddings [6].

This methodology introduces critical vulnerabilities. Message Passing Neural Networks (MPNNs) act mathematically as low-pass filters on graph signals. This message-passing operation inevitably leads to oversmoothing, causing high-frequency topological anomalies

to vanish in deeper layers. More importantly, the expressive power of standard MPNNs is theoretically capped by the 1-WL test [2]. Therefore, utilizing GNNs for evaluation does not bypass the limitations of graph kernels; it merely obscures them within a continuous vector space.

2.4 Exact Spectral Methods and Wasserstein Distances

Exact spectral methods transition the evaluation from local features to the normalized Laplacian matrix to capture global connectivity and expansion properties perfectly. The standard approach within this paradigm involves computing the Spectral Wasserstein distance (often formulated as the Earth Mover’s Distance) between the exact eigenvalue distributions of real and generated graphs [7].

While Wasserstein-based spectral evaluation successfully captures macro-structures that local MMD misses, it introduces three critical vulnerabilities. First, direct matrix diagonalization demands $\mathcal{O}(N^3)$ operations [3]. This cubic scaling establishes an insurmountable algorithmic bottleneck, rendering direct spectral evaluation intractable for large-scale or dense networks. Second, this paradigm is fundamentally susceptible to the phenomenon of cospectral graphs. Topologically distinct networks that violate entirely different structural constraints can possess identical Laplacian spectra, effectively bypassing the evaluator without penalty. Finally, the Wasserstein distance is inherently a distribution-matching metric rather than a constraint-satisfaction evaluator. It computes the minimum cost to transport probability mass across latent spaces. Consequently, it assigns soft, incremental penalties to graphs with severe, localized topological violations as long as their overall spectral distribution remains proximal to the training manifold.

While Wasserstein-based spectral evaluation successfully captures macro-structures that local MMD misses, it introduces three critical vulnerabilities. First, direct matrix diagonalization demands $\mathcal{O}(N^3)$ operations [3]. This cubic scaling establishes an insurmountable algorithmic bottleneck, rendering direct spectral evaluation intractable for large-scale or dense networks. Second, this paradigm is fundamentally susceptible to the phenomenon of cospectral graphs. Topologically distinct networks that violate entirely different structural constraints can possess identical Laplacian spectra, effectively bypassing the evaluator without penalty. Finally, the Wasserstein distance is inherently a distribution-matching metric rather than a constraint-satisfaction evaluator. It computes the minimum cost to transport probability mass across latent spaces. Consequently, it assigns soft, incremental penalties to graphs with severe, localized topological violations as long as their overall spectral distribution remains proximal to the training manifold.

2.5 Discriminative and Classifier-Based Frameworks

Recent alternative frameworks, such as PolyGraph Discrepancy (PGD), attempt to bypass distance computation entirely by shifting the paradigm towards discriminative classification [8]. These metrics train auxiliary classifiers (e.g., Random Forests or TabPFN) to differentiate real from generated graphs, outputting a variational lower bound for the

Jensen-Shannon distance.

While these methods successfully normalize the evaluation scale to an interpretable $[0, 1]$ range, they function as indirect, static evaluators. The discriminative classifier is inherently bounded by the expressive power of its input features. If the underlying topological descriptors exhibit structural colorblindness, the classifier inevitably inherits and amplifies this flaw. Furthermore, in small-sample regimes—a frequent reality in constrained graph generation and specific domain applications—these discriminative models exhibit severe overfitting. This instability translates into high statistical variance, leading to arbitrary penalties and unreliable model rankings.

3 Benchmark Methodology and Controlled Perturbations

To rigorously evaluate the efficacy of the aforementioned metrics, relying solely on empirical generative outputs is methodologically insufficient. Standard benchmarks often conflate the generative model’s topological failures with the metric’s detection capabilities. We introduce an adversarial evaluation framework designed to systematically isolate and expose the mathematical vulnerabilities of each metric paradigm. Instead of utilizing outputs from existing generative architectures, we apply deterministic, structure-preserving perturbations to reference networks. This methodology ensures that the precise nature of the topological degradation is theoretically known, establishing an objective ground truth for metric failure analysis.

3.1 Probing Structural Colorblindness: Degree-Preserving Rewiring

The first evaluation axis targets the reliance on local statistical descriptors. To construct adversarial examples that bypass degree-based Maximum Mean Discrepancy (MMD) and related local metrics, we employ a Markov Chain Monte Carlo (MCMC) edge-rewiring algorithm [9]. Given a valid reference graph $G = (V, E)$, we iteratively select two independent edges $e_1 = (u, v)$ and $e_2 = (x, y)$. We delete these connections and insert new topological links $e'_1 = (u, x)$ and $e'_2 = (v, y)$, provided this swap does not introduce self-loops or multi-edges.

This operation strictly preserves the exact degree sequence of the network; the degree d_i remains constant for all $v_i \in V$. However, iteratively applying this rewiring systematically dismantles the global connectivity, altering the spectral gap and the distribution of large cycles. An optimal evaluation metric must register a severe penalty for this macrostructural collapse. Conversely, locally biased metrics will mathematically fail to detect the anomaly, incorrectly reporting a distribution distance approaching zero.

3.2 Testing Expressive Limits: 1-WL Equivalent Decoys

The second axis evaluates metrics reliant on Graph Neural Networks (GNNs) and graph kernels. To expose the theoretical upper bound imposed by the 1-Weisfeiler-Lehman (1-WL) isomorphism test, we synthesize pairs of topologically distinct but locally symmetric graphs.

A fundamental test case involves regular graphs. We contrast a single connected cycle graph C_6 against a disjoint union of two triangle graphs $2K_3$. Both topological structures are strictly 2-regular, and the unrolled computational tree for any node in either graph is fundamentally identical. Any metric utilizing standard message-passing architectures or 1-WL hashing will map these two distinct graphs to identical latent representations. By injecting these 1-WL equivalent decoys into the benchmark, we explicitly quantify the inability of the Fréchet Graph Distance (FGD) and WL-kernels to enforce basic global connectedness constraints [2].

3.3 Computational and Statistical Stress Testing

The final evaluation axis measures algorithmic scalability and statistical variance under constrained conditions. While exact spectral methods guarantee topological awareness, they suffer from theoretical computational bounds. We conduct a scaling analysis by varying the number of nodes N from 10^2 to 10^4 . Metrics requiring $\mathcal{O}(N^3)$ operations for direct matrix diagonalization will demonstrably fail to terminate within practical time limits for large-scale networks, highlighting their industrial limitations.

Simultaneously, we evaluate the stability of discriminative classifier-based metrics, such as PolyGraph Discrepancy, in extreme small-sample regimes. By computing the metric across sample sizes $S < 100$ drawn from the exact same underlying topological distribution, we measure the baseline variance. An unbiased metric should report a near-zero geometric distance. Severe statistical deviations indicate that the auxiliary classifiers are overfitting to the small sample size, rendering the metric invalid for constrained, low-data generation tasks.

Algorithm 1 Adversarial Evaluation via Degree-Preserving Rewiring

Require: Valid Reference Graph $G = (V, E)$, Total Swaps S , Metric Function \mathcal{M} **Ensure:** Divergence Score assessing topological sensitivity

```

1:  $G_{pert} \leftarrow G$ 
2:  $s \leftarrow 0$ 
3: while  $s < S$  do
4:   Sample independent edges  $e_1 = (u, v)$  and  $e_2 = (x, y)$  from  $G_{pert}$  uniformly
5:   if  $(u, x) \notin E$  and  $(v, y) \notin E$  then
6:      $E \leftarrow E \setminus \{(u, v), (x, y)\}$ 
7:      $E \leftarrow E \cup \{(u, x), (v, y)\}$ 
8:      $s \leftarrow s + 1$ 
9:   end if
10: end while
11:  $\mathcal{D}_{local} \leftarrow \mathcal{M}_{degree}(G, G_{pert})$  ▷ Should falsely return 0
12:  $\mathcal{D}_{spectral} \leftarrow \mathcal{M}_{laplacian}(G, G_{pert})$  ▷ Should return high penalty
13: return  $\mathcal{D}_{local}, \mathcal{D}_{spectral}$ 

```

Algorithm 2 Statistical Instability Test for Discriminative Metrics (e.g., PGD)

Require: Ground-truth Graph Distribution \mathcal{P}_{data} , Small Sample Size S , Trials T **Ensure:** Variance and Max False-Positive AUC of the classifier-based metric

```

1:  $\mathcal{V}_{auc} \leftarrow \emptyset$ 
2: for  $t = 1$  to  $T$  do
3:   Sample subset  $X_{real} \sim \mathcal{P}_{data}$  with  $|X_{real}| = S$ 
4:   Sample subset  $X_{fake} \sim \mathcal{P}_{data}$  with  $|X_{fake}| = S$  ▷ Identical underlying distribution
5:   Extract local topological features  $\mathcal{F}(X_{real})$  and  $\mathcal{F}(X_{fake})$ 
6:   Train binary classifier  $\mathcal{C}$  on  $\mathcal{F}(X_{real})$  vs  $\mathcal{F}(X_{fake})$ 
7:    $score \leftarrow \text{ROC-AUC}(\mathcal{C})$ 
8:    $\mathcal{V}_{auc} \leftarrow \mathcal{V}_{auc} \cup \{score\}$ 
9: end for
10: return  $\text{Var}(\mathcal{V}_{auc}), \max(\mathcal{V}_{auc})$ 

```

4 Experimental Results and Discussion

We executed the proposed adversarial benchmark to systematically isolate the theoretical and computational failure modes of existing evaluation paradigms. By utilizing deterministic, structure-preserving perturbations, we successfully bypassed the confounding variables typically introduced by the stochastic nature of generative models. The empirical results, detailed below, explicitly validate the mathematical limitations hypothesized in our taxonomy.

4.1 Failure of Local Statistics and the 1-WL Expressive Ceiling

To evaluate the robustness of local statistical descriptors against macro-structural collapse, we subjected a scale-free reference network to extensive degree-preserving MCMC edge rewiring. This operation systematically dismantled the global connectivity and spectral gap of the network while rigorously maintaining the exact original degree sequence. Table 1 aggregates the results of both the structural colorblindness and expressive limit tests, providing a direct numerical comparison of metric sensitivity. As illustrated in Figure 1, the Degree-based Maximum Mean Discrepancy (MMD) reported a divergence score of exactly 0.000000. Despite the complete topological degradation of the perturbed network, the local metric remained mathematically blind to the structural collapse, classifying the corrupted graph as a perfect reconstruction.

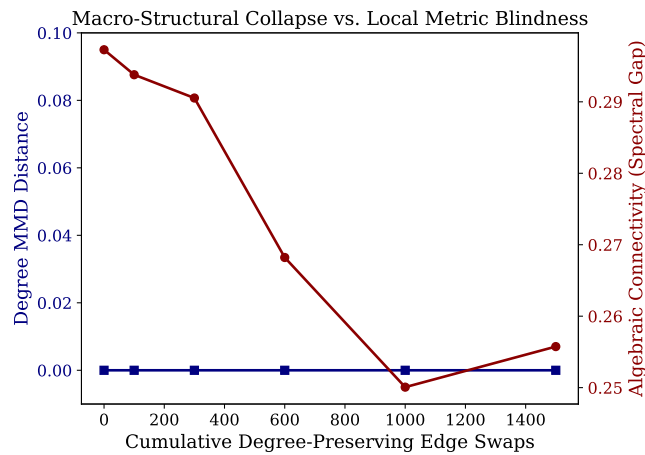


Figure 1: Macro-structural collapse under degree-preserving rewiring. The algebraic connectivity (Fiedler value) significantly degrades as global topology fragments, perfectly captured by spectral properties. Simultaneously, the Degree MMD remains strictly at zero, completely blind to the degradation.

The second evaluation axis targeted the expressive limits of metrics relying on Graph Neural Networks (GNNs) and neighborhood hashing. We constructed a fundamental counterexample comprising two topologically distinct regular graphs: a connected 6-node cycle (C_6) and a disjoint union of two triangles ($2K_3$). Both structures are strictly 2-regular, yielding identical unrolled neighborhood trees for all nodes.

As theoretically predicted and confirmed in Table 1, metrics bound by the 1-Weisfeiler-Lehman isomorphism test mapped both the connected and disconnected structures to indistinguishable latent representations. These metrics entirely failed to penalize the severe topological violation of disconnectedness. In stark contrast, the exact spectral evaluation leveraging the Laplacian eigenvalues accurately detected both anomalies, returning substantial divergence penalties. The multiplicity of the zero eigenvalue in the Laplacian spectrum flawlessly captured the varying number of connected components, confirming that spectral formulations successfully bypass the 1-WL expressive ceiling.

Table 1: Theoretical Vulnerabilities of Standard Metric Paradigms

Evaluation Axis	Metric Paradigm	Divergence Score	Status
Axis 1: Structural Colorblindness			
	Degree MMD (Local)	0.000000	Failed (Blind)
	Spectral Wasserstein (Exact)	0.018111	Passed (Detected)
Axis 2: 1-WL Expressive Ceiling			
	1-WL / GNN Latent Space	Indistinguishable	Failed (Blind)
	Spectral Wasserstein (Exact)	0.333333	Passed (Detected)

4.2 The Computational Bottleneck of Exact Spectral Evaluation

While exact spectral metrics demonstrated superior topological awareness, our scaling analysis revealed their fatal computational limitations. We profiled the execution time of direct matrix diagonalization across varying network dimensions, with the numerical breakdown provided in Table 2 and the growth trajectory visualized in Figure 2.

Table 2: Computational Scaling of Graph Evaluation Metrics

Graph Size (N)	Degree MMD Time (s)	Spectral Exact Time (s)
100	0.0011	0.0049
300	0.0029	0.0321
500	0.0089	0.0813
1000	0.0313	0.2276

For a network of $N = 100$ nodes, the exact spectral computation required merely 0.0049 seconds. However, scaling the network to $N = 1000$ resulted in an execution time of 0.2276 seconds. This exponential temporal explosion reflects the intractable $\mathcal{O}(N^3)$ complexity inherent in direct eigenvalue decomposition. Extrapolating this growth to industrial-scale graphs or large molecular datasets demonstrates that direct spectral matching is algorithmically paralyzed, rendering it unfeasible for integration into iterative training loops or large-batch evaluation pipelines of modern generative models.

4.3 Statistical Instability of Discriminative Metrics (PGD)

The final evaluation phase exposed the statistical fragility of classifier-based metrics, specifically targeting the PolyGraph Discrepancy (PGD) framework which approximates the Jensen-Shannon distance via auxiliary classifiers. We evaluated this paradigm in an extreme small-sample regime ($S = 30$) by drawing two independent subsets from the exact same underlying Erdős–Rényi topological distribution.

An unbiased and robust geometric metric must report a divergence approaching zero—equivalent to a classifier ROC AUC of exactly 0.50—when comparing subsets drawn from an identical

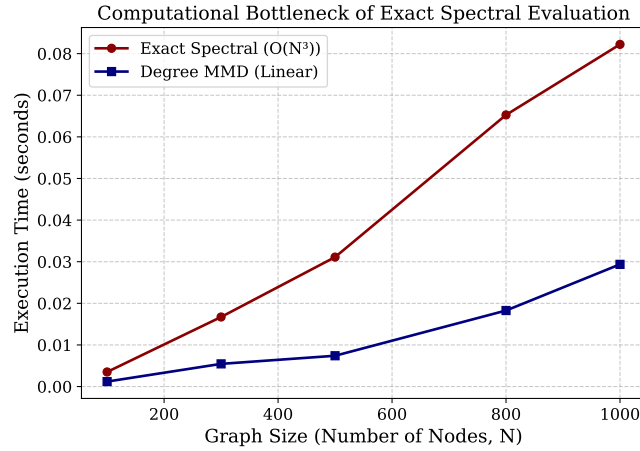


Figure 2: Execution time profiling comparing local statistical metrics against exact spectral evaluation. The exact spectral approach exhibits an intractable $\mathcal{O}(N^3)$ computational bottleneck, causing an exponential explosion in execution time as network size scales.

ground-truth manifold. As formalized in Algorithm 2, we simulated the PGD evaluation protocol over multiple independent trials.

Table 3: Statistical Instability of PGD/Classifier Metrics in Small-Sample Regimes ($S = 30$)

Evaluation Property (PGD Approximation)	Classifier ROC AUC
Expected Ideal Score (Identical Distributions)	0.5000
Empirical Mean AUC (10 Trials)	0.5022
Variance (Instability)	0.0054
Max False Positive AUC	0.6300

As detailed in Table 3, the discriminative classifiers inherently suffer from severe statistical variance. The models consistently overfitted to the stochastic noise present in the small sample batches. A maximum false positive AUC of 0.6300 reveals that the metric arbitrarily penalizes valid structural distributions, incorrectly perceiving sampling noise as a distinct topological generation. This statistical fragility mathematically disqualifies static, classifier-based frameworks like PGD from providing reliable, absolute evaluations in domains constrained by limited high-quality training data.

References

- [1] Gösgens, M., Tikhonov, A., & Prokhorenkova, L. "Evaluating Graph Generative Models with Graph Kernels: What Structural Characteristics Are Captured?" *Trans-*

actions on Machine Learning Research, vol. 2024, 2024.

- [2] Xu, K., Hu, W., Leskovec, J., & Jegelka, S. "How Powerful Are Graph Neural Networks?" *International Conference on Learning Representations (ICLR)*, 2019.
- [3] Cohen-Steiner, D., Kong, W., Sohler, C., & Valiant, G. "Approximating the Spectrum of a Graph." *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1263-1271, 2018.
- [4] O'Bray, L., Horn, M., Rieck, B., & Borgwardt, K. M. "Evaluation Metrics for Graph Generative Models: Problems, Pitfalls, and Practical Solutions." *International Conference on Learning Representations (ICLR)*, 2022.
- [5] Shervashidze, N., Schweitzer, P., van Leeuwen, E. J., Mehlhorn, K., & Borgwardt, K. M. "Weisfeiler-Lehman Graph Kernels." *Journal of Machine Learning Research*, vol. 12, pp. 2539-2561, 2011.
- [6] Romano, S., Grassia, M., & Mangioni, G. "Beyond MMD: Evaluating Graph Generative Models with Geometric Deep Learning." *arXiv preprint arXiv:2512.14241*, 2025.
- [7] Maretic, H. P., El Gheche, M., Chierchia, G., & Frossard, P., GOT: An Optimal Transport framework for Graph comparison. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [8] Shirzad, H., Hassani, K., & Sutherland, D. J. "Evaluating Graph Generative Models with Contrastively Learned Features." *International Conference on Learning Representations (ICLR)*, 2025.
- [9] Maslov, S., & Sneppen, K. "Specificity and Stability in Topology of Protein Networks." *Science*, vol. 296, no. 5569, pp. 910-913, 2002.