

Oversmoothing, “Oversquashing”, Heterophily, Long-Range, and more: Demystifying Common Beliefs in Graph Machine Learning

Adrian Arnaiz-Rodriguez¹ and Federico Errica²

¹ ELLIS Alicante – adrian@ellisalicante.org

² NEC Laboratories Europe – federico.errica@necclab.eu

Abstract. After a renaissance phase in which researchers revisited the message-passing paradigm through the lens of deep learning, the graph machine learning community shifted its attention towards a deeper and practical understanding of message-passing’s benefits and limitations. In this **position** paper, we notice how the fast pace of progress around the topics of oversmoothing and oversquashing, the homophily-heterophily dichotomy, and long-range tasks, came with the consolidation of commonly accepted beliefs and assumptions that are not always true nor easy to distinguish from each other. We argue that this has led to ambiguities around the investigated problems, preventing researchers from focusing on and addressing precise research questions while causing a good amount of misunderstandings. Our contribution wants to make such common beliefs explicit and encourage critical thinking around these topics, supported by simple but noteworthy counterexamples. The hope is to clarify the distinction between the different issues and promote separate but intertwined research directions to address them.

Keywords: oversmoothing · oversquashing · heterophily · long-range propagation · graph neural networks · graph machine learning

1 Introduction

The last decade has seen an increasing scholarly interest in machine learning for graph-structured data [10, 44, 73, 92]. After an initial focus on the design of various message-passing architectures [39], inheriting from the recurrent [87] and convolutional [72] Deep Graph Networks (DGN), together with the analysis of their expressive power [74, 103], researchers later turned their attention to the intrinsic limitations of the message-passing strategy and the relation between the graph, the task, and the attainable performance. We refer, in particular, to the fact that node embeddings may become increasingly similar to each other as more message-passing layers are used [85], the loss of information that results from aggregating too many messages onto a single node embedding [4], the presence of topological bottlenecks [95], the existence of neighbors of different classes [99], and the propagation of information between far ends of a graph [31]. Addressing these limitations makes a difference when applying message-passing models, including

foundational ones [14], to large and topologically varying graphs at different scales, from proteins with hundreds of thousands of atoms [108] to dynamic social networks [64] with billions of users, where such limitations manifest together.

A pace of research so rapid can sometimes lead, however, to the premature consolidation of ideas and beliefs that have not been thoroughly verified. There are several reasons for this to happen: the (perhaps too) intense pressure to publish, follow the latest scientific trends, and demonstrate state-of-the-art performance. As a result, we may end up putting the spotlight on positive findings but overlooking contradictory evidence, eventually accepting hypotheses as canon. In this **position paper**, we argue and provide evidence that this is potentially the case for the afore-mentioned issues of *oversmoothing* (OSM), *oversquashing* (OSQ), *heterophily*, and *long-range dependencies*, driving researchers away by the difficulty of circumventing common beliefs while introducing novel contributions. Scientific progress is therefore slowed down both in terms of reduced workforce and clarity of the problems to be addressed; failure to acknowledge existing inconsistencies may well lead to a reiterated spreading of questionable claims.

While reviewing the literature, we identify and isolate nine common beliefs that, in our opinion, cause great confusion and ambiguities in the field. We then demystify such beliefs by providing *simple* and possibly memorable counterexamples that should be easy to recall. By encouraging critical thinking around these issues and separating the different research questions, we hope to foster further advancements in the graph machine learning field.

Table 1 summarizes our findings about common beliefs in the literature, together with the list of papers where we could find mentions of them. We logically divide common beliefs about OSM, OSQ, and the homophily-heterophily dichotomy. In the following sections, we discuss each belief, provide counterexamples, and summarize our arguments with take-home messages.

2 Is Oversmoothing Really a Problem?

Oversmoothing broadly refers to the phenomenon where, as we stack more message-passing layers in a DGN, the node embeddings become increasingly similar to each other eventually collapsing into a low-dimensional [51, 77]–or even single-vector [19]–subspace. This creates an almost constant representation, independent of the original node-feature distribution, and can **potentially** result in loss of discriminative power along the way [19, 59, 77].

Formally, let $H^\ell \in \mathbb{R}^{n \times d}$ be the matrix of node embeddings after ℓ message-passing layers in a DGN, where n is the number of nodes and d is the hidden dimension. Consider a similarity (or separation) function $\pi: \mathcal{H} \rightarrow \mathbb{R}$, where $\mathcal{H} = \mathbb{R}^{n \times d}$ is the space of all possible embedding matrices. We say that a DGN experiences OSM if

$$\lim_{\ell \rightarrow \infty} \pi(H^\ell) = c. \quad (1)$$

where c is some constant indicating a collapse of embeddings. Deviation metrics (e.g., Dirichlet Energy [19], MAD [22]) and subspace-collapse criteria [50, 77, 84]

Table 1: List of common beliefs together with papers that make those claims.

	Common Belief	References
OSM	1. OSM is the cause of performance degradation.	[17–19, 22, 24, 25, 30, 33, 47, 48, 51–54, 59, 61, 70, 76, 83–86, 93, 97, 98, 100, 101, 109, 111, 112]
	2. OSM is a property of all DGNs	[1, 3–5, 7–9, 18, 24, 25, 29, 30, 35, 36, 38, 42, 47, 48, 51–53, 55, 58, 76, 84–86, 89, 95, 100–102, 104, 109, 111, 112]
Hom-Het	3. Homophily is good, heterophily is bad.	[5, 8, 9, 15, 17, 40, 43, 60, 63, 66, 79, 80, 84, 98, 110, 111, 113]
	4. Long-range propagation is evaluated on heterophilic graphs.	[3, 5, 8, 16, 42, 50, 70, 96, 97]
	5. Different classes imply different features.	[2, 8, 9, 15, 50, 60, 66, 67, 69, 79, 85, 94, 99, 110]
OSQ	6. OSQ synonym of a topological bottleneck.	[3, 5, 7–9, 11–13, 16, 25, 28–30, 35–37, 41, 42, 45, 50, 53, 54, 62, 76, 89–91, 93–97, 104, 106]
	7. OSQ synonym of computational bottleneck.	[1, 3–5, 7–9, 11–13, 16, 25, 28, 30, 31, 35, 42, 45, 46, 50, 54, 76, 89–91, 93–95, 97]
	8. OSQ problematic for long-range tasks.	[1, 3, 4, 7, 9, 12, 13, 16, 18, 25, 28, 29, 35–37, 41, 45, 50, 54, 62, 76, 90, 91, 93, 95, 97, 106]
	9. Topological bottlenecks associated with long-range problems.	[3, 16, 36, 54, 62, 90, 95, 97]

all measure the same OSM intuition: as depth increases, node embeddings shrink toward a nearly constant subspace or degree-proportional vectors.

2.1 Belief: OSM is a Property of All DGNs.

A widespread claim in the literature is that *OSM happens regardless of the specific architecture or the underlying graph*. Early theoretical works support this view by analyzing message-passing propagation as a diffusion process: iterated normalized-Laplacian updates converge to a degree-weighted stationary distribution [19, 42], while heat-kernel diffusion converges to a constant vector [6, 77]. The resulting bounds quantify the rate of OSM in terms of the singular values of the feature transform W and the eigenvalues of the graph structure G .

These conclusions, however, rely on restrictive assumptions. Later work has relaxed them by introducing learnable feature transforms, non-linear activations, and more elaborate architectures [19, 77, 101], yet no existing proof shows inevitable collapse under realistic training regimes. In practice, remedies such as residual/skip connections, normalization layers, or gating mechanisms are explicitly architectural changes designed to maintain local distinctions, calling into question the universality of this OSM claim.

In addition, many studies probe OSM with *untrained* (weights frozen at initialization) linear GCN stacks, which is an experimental choice that hides the effect of learning and may lead to the wrong conclusions, as noted by [107]. Indeed, [27, Figure 2] report OSM only for frozen-weight networks; once parameters are allowed to adapt, the models preserve informative variance.

Together, these observations suggest that OSM is not an inevitable consequence of message-passing but rather a contingent outcome that depends on training dynamics and architectural design choices.

Empirical Example We show a simple training scenario where we see how OSM is *not* a property of all DGNs and how different elements make it difficult to draw clear conclusions. We train several DGNs under two propagation variants: the vanilla AXW update and the rescaled $AX(2W)$, inspired by [84, Figure 1]. We measure OSM with 2 different metrics: Dirichlet Energy (DE) and its norm-normalized version, the Rayleigh Coefficient (RQ), which was also previously used in some works [18, 40, 70, 84]. Figure 1 shows that: (i) some architectures never collapse, (ii) a minor rescaling can reverse the trend, and (iii) DE and RQ often disagree. Hence OSM is neither universal nor straightforward to diagnose.

First, OSM, as measured by Dirichlet Energy (DE), is not universal: GIN’s DE explodes instead of collapsing in the vanilla setting (a), which shows that changing the aggregation-function may lead to an opposite behavior of OSM effect. Second, a minor scaling in the feature transformation ($2W$ instead of W) flips the behavior of several models: curves that decayed in (a) now grow or stabilize, and vice-versa. Similar small tweaks (normalization layers, self-loops, alternative aggregators) can likewise create or remove DE collapse, which has been leveraged by prior work to propose OSM mitigation approaches, such as the ones based on feature normalization [109, 111].

Finally, whether OSM is observed depends on the measure of choice. DE reflects raw smoothness, whereas the RQ normalizes by the feature norm. As a result, the same model can exhibit mutually contradictory trends for different metrics [107]. First, GCN’s DE collapses under the vanilla aggregation (a), explodes with a simple rescaling (b), yet RQ remains essentially flat in both normalised plots (c–d), indicating that GCN embeddings are being rescaled, not necessarily oversmoothed. On the contrary, GAT and SAGE, which had similar behavior as GCN in (a) and (b), decay using RQ (c-d) at a (roughly) linear rate, highlighting how different architectures respond differently to the use of RQ instead of DE. Furthermore, subfigure (d) shows that $AX(2W)$ can stabilize RQ for certain methods, suggesting that normalizations or small parameter adjustments do not affect all models uniformly. Therefore, conclusions about OSM depend on metric *and* model, a rare observation in the literature.

In conclusion, OSM is neither inevitable nor uniquely defined: its observation hinges on the architecture, on seemingly innocuous hyper-parameters, and on which stability metric (DE vs. RQ) one adopts. Therefore, a natural question arises: does OSM actually limit the models’ predictive accuracy? We investigate this question in the next section.

2.2 Belief: OSM is the Cause of Performance Degradation.

Part of the literature focuses on the narrative that OSM is the cause of lower test accuracy in DGNs. The hypothesis is that if embeddings collapse to a non-meaningful space, then the separability of the nodes will become challenging and accuracy decreases.

However, this hypothesis ignores some critical aspects, such as *i)* the separability of node embeddings with respect to the nodes’ labels, and *ii)* how

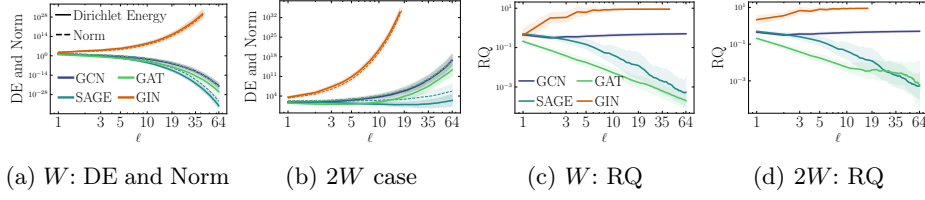


Fig.1: **(a-b)**: We depict the evolution, with increasing number of layers, of the $DE = \text{tr}(X^T \Delta X)$ and the feature norm $\|X\|_F^2$, using W and $2W$ feature transformations for different architectures. **(c-d)**: Evolution of the $RQ = \text{tr}(X^T \Delta X) / \|X\|_F^2$ for W and $2W$ as before. Experiments run on the Cora dataset for 50 random seeds. A larger version for better visualization is available in Fig. 7.

such separability evolves in the intermediate OSM phase (*if* it happens, as we discussed before).

Regarding the first statement, although it is true that if there is total collapse to the same value then the embeddings will not be informative at all, the main problem remains the node embeddings’ separability. As shown in the previous section, some changes in the architectures can avoid OSM, but they might have no impact on the overall accuracy. For instance, multiplying by two the weight matrix leads to a general increase in DE for all architectures, however, the accuracy will remain the same since the embeddings have been simply scaled up and the embeddings’ separability is not affected negatively. In addition, avoiding an embedding collapse does not necessarily lead to an improvement in generalization accuracy. For instance, comparing a GCN with and without bias, both versions show a decrease in performance as the number of layers increases, whereas only GCN shows a collapse in DE [85, Figure 3].

On the other hand, and related to the second statement, embedding collapse will not always lead to a decrease in accuracy. OSM happens faster in some subspaces than in others and this effect will be beneficial if labels are correlated with those subspaces [55]. For instance, if we classify points into two classes, and all nodes of distinct classes collapse into different points, the OSM metric will detect such collapse. However, the separability of the node embeddings will remain possible, illustrating also the limits of wide-spread used OSM metrics with respect to label information.

This intuitive behavior has been identified in the literature as a form of “beneficial” smoothing phase [55, 84, 102]). In this phase, the nodes of each class first collapse into a class-dependent point before the *potential* second stage, at which point all nodes converge to the same representation. Finally, although overall pairwise distances might shrink in deep layers, *within-class* distances might contract more than *between-class* ones, so class separability improves despite the global collapse, as discussed by [27].

In conclusion, low accuracy in DGNs cannot be attributed to OSM alone. The separability of node embeddings plays a major role, where other training

problems such as vanishing gradients or over-fitting also arise when using a big number of layers [7, 27, 78, 105, 109].

Message of the Section

- i) OSM is not a property of all DGNs
- ii) OSM is not necessarily the cause of performance degradation. The performance is related to node embeddings’ separability, which can be also affected by many other elements, such as vanishing gradients.
- iii) Therefore, to study the performance of DGNs, it might be better to study how they achieve separability of node embeddings, and how the OSM relates to node separability.

3 Homophily-Heterophily and the Role of the Task

In the context of node classification, the term *homophily* (resp. *heterophily*) generally refers to some form of similarity (resp. dissimilarity) between a node and its neighbors [71]. This (dis)similarity can be measured with respect to class labels, node features, or both; the vast majority of works in the literature opt for the first, *but this choice is often implicit and taken for granted*, making some statements hard to interpret when one is aware of the other ways to measure it.

3.1 Belief: Homophily is Good, Heterophily is Bad

A recurrent narrative in the literature is that the message-passing mechanism of DGNs is particularly suited for homophilic graphs, whereas it is unfit for heterophilic graphs. The apparent motivation is that, in homophilic graphs, all you need to do to solve a node classification task is to look at similar neighbors, and a local message-passing strategy implements just the right inductive bias. This is in contrast to a class-heterophilic graph, where there exist neighboring nodes of a different class that might make it harder for message-passing to isolate the “relevant information”, intended as neighbors of the same class. Such a belief is supported by empirical evidence on a rather restricted set of benchmarks [75, 79, 88] with varying levels of homophily/heterophily.

Researchers already tried to challenge these considerations in the past [34, 68, 69, 81, 99]. Consider Figure 2 (left), where a bipartite graph with identical node features is fully class-heterophilic. If we apply a single sum-based graph convolution, the nodes can be perfectly classified as the resulting embeddings depend on the incoming degree. Therefore, **there exist heterophilic graphs where a DGN can achieve perfect classification**. On the contrary, Figure 2 (right) depicts a highly class-homophilic graph, where nodes belong to one of two classes if they are at a distance greater or lower than five from a given node s . In this case, the information on the nodes does not even matter; if we were to follow the above belief, we would be encouraged to use a few layers of message-passing, and as a result **we could never solve the task perfectly**.

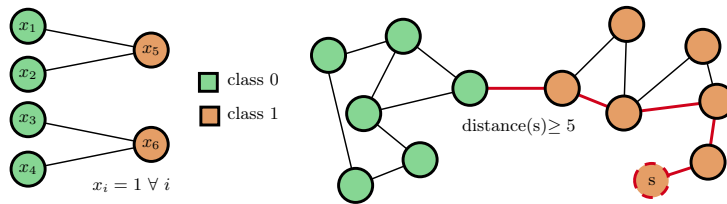


Fig. 2: *Left*: a fully heterophilic graph inspired by [69] where a 1-layer, sum-based DGN can perfectly classify the nodes due to a difference in the node degree. *Right*: a highly homophilic graph where the task is to predict if a node is at a distance greater than five from a specific node. Here, the performances of a DGN will be poor unless information from nodes of another community – from the perspective of a class-0 node – is captured.

3.2 Belief: Different Classes Imply Different Features

In the previous section, we deliberately ignored the interplay between a node’s features and its class label, which is induced by the task at hand. The reason is that we wanted to clarify the distinction with another, more subtle, and problematic belief: nodes belonging from different classes should have different (to be read as separable) feature distributions. Under this assumption, it is also implicit that class homophily will imply feature homophily.

Such an assumption is often key in arguments supporting the belief of Section 3.1. Indeed, if nodes of different classes have different feature distributions, then applying a local message-passing iteration to a highly class-homophilic graph should “preserve the distance” between node embeddings of different classes. On the contrary, in a heterophilic setting, a convolution would mix information coming from different feature distributions, which may be detrimental to performances.

The logic is not incorrect per-se, but our key counterargument is the following: if different classes imply different feature distributions, why do we need a DGN rather than a simple MLP? **Either there is a very strong assumption** that the task does not depend on the topological information, or the feature-class distributions induced by the task allow us to somehow take shortcuts, neglecting the role that the topology might have. Yet, it is important to mention the denoising/regularizing effects of DGNs in semi-supervised scenarios [34, 49].

It appears therefore necessary to consider less trivial and more fine-grained scenarios, where the feature distributions of different classes partially or totally overlap, the topology has a key role in the task definition, and topological properties induce a positive/negative effect on the performance of message-passing models as done, for instance, in recent works [20, 110].

3.3 Belief: Long-range Propagation is Evaluated on Heterophilic Graphs

The common beliefs of Sections 3.1 and 3.2 have been used to support yet another argument, namely that we should evaluate the ability of DGNs to propagate

long-range information on heterophilic graphs. The rationale seems to be that, in order for DGNs to perform well, nodes of a given class should focus on information of similar (w.r.t. class and/or features) nodes; therefore, in a heterophilic graph, it may be necessary to capture information far away (i.e., long-range) from the immediate neighborhood.

Once more, what is really important is to **distinguish the task**, e.g., one that depends on long-range propagation, **from the class labels the task induces on the nodes**. As a matter of fact, the “long-range” task of Figure 2 (right) induces a highly homophilic graph, while the heterophilic graph of Figure 2 (left) is not associated with a long-range propagation task. Therefore, we cannot draw a relation between long-range tasks and heterophily without further assumptions.

Message of the Section

- i) Generic claims about the performance of DGNs under homophily and heterophily do not hold, nor does their relation with long-range problems.
- ii) Homophily/heterophily is a function of the task, but the converse is not true.
- iii) We should move past the coarse homophily-heterophily dichotomy and **focus more on the task** and the interplay it induces between features, structure, and class labels.

4 The Many Facets of “Oversquashing” and their Negative Implications

The term oversquashing originated from [4] and referred to an “*exponentially growing information into fixed-size vector*” by repeated application of message-passing. In other words, oversquashing was associated with the **computational tree** (Figure 3) induced by message-passing layers on each node of the graph. Later, oversquashing was connected by [95] to the existence of **topological bottlenecks**: “*edges with high negative curvature are those causing the graph bottleneck and thus leading to the over-squashing phenomenon*”. Since then, researchers have adopted one or even both definitions of oversquashing at the same time, contributing to an apparent understanding that these definitions subsume the same problem. In this section, we argue that **this is not the case** and that, as a community, **we should clearly separate the term “oversquashing”** into (at least) two separate terms:

Computational Bottlenecks and Topological Bottlenecks

Computational bottlenecks, defined in 1, are inevitably related to the *message-passing architecture*, for which the graph is the computational medium, whereas topological bottlenecks refer to the *graph connectivity*. These two problems are clearly intertwined, but in the following we show that they do not always coexist, hence it makes sense to treat them as **fundamentally distinct problems**.

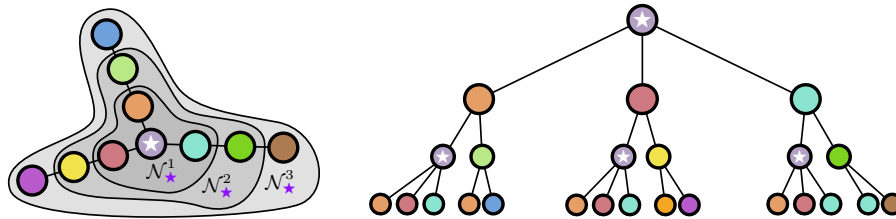


Fig. 3: We intuitively visualize what happens when we repeatedly aggregate the neighborhood of the star node using the message-passing paradigm. Appendix A.4 defines a computational bottleneck as the size of this computational tree (24 computational nodes *vs* nine 3-hop neighbors).

4.1 Belief: Oversquashing as Synonym of Topological Bottleneck

The prolific line of work that associates “oversquashing” with topological bottlenecks seems to have gained popularity with [95]. In that paper, edges with negative curvature are first associated with (topological) bottlenecks, then a theorem puts in relation message-passing on a graph containing a bottleneck with the Jacobian sensitivity of node representations, typically defined by $I(u, v) = \|\partial \mathbf{h}_u^K / \partial \mathbf{h}_v^0\|$ and denoting how much the final representation of a node u after K layers is influenced by the initial representation of a node v [104]. Put simply, a topological bottleneck may imply low sensitivity. **However, the converse is not necessarily true:** we can have low sensitivity on a graph where there are no bottlenecks. Figure 4 (left) shows a grid graph where there are no topological bottlenecks. The repeated application of message passing will, however, quickly generate a computational bottleneck. Therefore, saying that there are no topological bottlenecks does not imply that there are no computational bottlenecks.

To improve on topological bottlenecks, a widely investigated approach is graph rewiring [9], which was also the subject of scrutiny recently [96, 97]. Rewiring is based on the intuition that improving topological bottlenecks metrics should improve the performance of DGNs [5, 12, 28, 54], by reducing the distance between nodes that should communicate. At the same time, it should become clear now that, under the DGN paradigm, *rewiring might worsen the computational bottleneck* – as long as the same number of message-passing layers is used – while improving the topological one. This perspective was also put forward by [35], with a theoretical analysis on how message filtering, as shown in Figure 2 (middle), reduces both the computational bottleneck and sensitivity yet improves performances while leaving the graph structure unaltered.

Another, slightly more technical way to see why low Jacobian sensitivity does not imply the presence of any topological bottlenecks is to follow the chain of upper bounds that link the metrics used to measure the computational and topological bottlenecks [5, 16, 40, 54]. Several recent works [16, 29] have shown that sensitivity is bounded by above by a term that includes the effective resistance, a purely topological distance metric that quantifies the expected commute time of

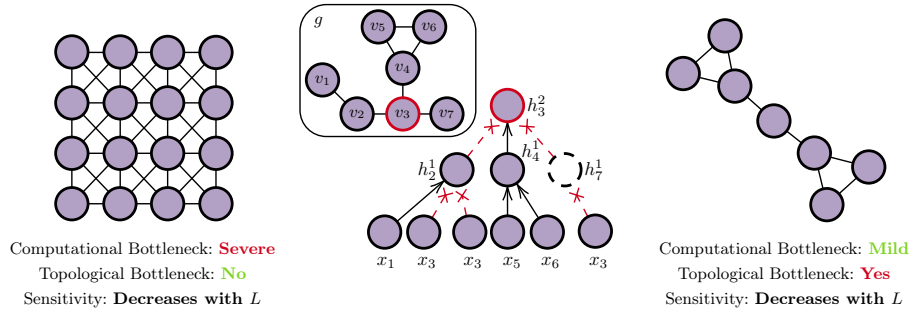


Fig. 4: *Left*: in a grid graph, the computational bottleneck grows very quickly but there is no topological bottleneck. *Middle*: A visualization of the computational graph rooted at node $v_3 \in \mathcal{V}_g$ for two message passing layers, highlighting how pruning messages reduces the computational bottleneck. *Right*: In this graph, there is a topological bottleneck and a mild computational bottleneck. As with the grid graph (Appendix B), the sensitivity decreases with the number of layers.

a random walk between nodes u and v [57]. In particular, the bound subtracts the lower bound on the maximum effective resistance, which depends on the inverse of the spectral gap (a proxy for topological bottlenecks) [21, 65]. This connection provides a useful intuition: as the topological bottleneck gets worse, the lower bound on the maximum effective resistance increases, which in turn reduces the upper bounds on the sensitivity of [16]. **However, the converse does not hold.** There are graphs where the maximum effective resistance between distant nodes is large despite the absence of topological bottlenecks. Consider the example of a grid graph in Figure 4 (left) where there is no topological bottleneck, yet the effective resistance between diagonally opposite corners grows linearly with the grid size. As a result, sensitivity between those nodes decays with depth, even though there are no topological bottlenecks. This illustrates that computational bottlenecks can arise independently of topological ones, and that low sensitivity **does not necessarily imply** the presence of either of them.

4.2 Belief: Oversquashing as Synonym of Computational Bottleneck

We briefly complement the previous section with a discussion on “oversquashing” as a computational bottleneck, which was introduced by [4] and has been the (often implicit) study subject of works that prune, to some extent, the computational tree induced on every node by the iterative message-passing process [35, 83]. Also in this case, there exist cases where reducing the computational bottleneck may be harmful: Figure 4 (right) shows a graph where there is a topological bottleneck but no severe computational bottleneck (for a limited number of layers). In this case, excessive pruning of the computational tree might cause distant nodes to interrupt all communications. Therefore, computational and topological bottlenecks are problems that should be tackled separately.

4.3 Belief: Oversquashing is Problematic for Long-range Tasks

Since its definition by [4], oversquashing has often been considered a problem in long-range tasks. The reason stems from its relation to the exponentially growing computational tree as the number of message-passing layers increases: whenever a node has to receive information from another node at distance d , classical (synchronous) message-passing architectures need to apply at least d layers to capture that information. As a result, the relevant information may get lost due to the exponentially large computational bottleneck. Importantly, topological bottlenecks can only make the problem worse, by forcing the information of a group of messages to be squeezed through an edge – see the next section for a discussion about long-range tasks and topological bottlenecks.

The main message here is that long-range tasks “force” classical message-passing architecture to create a computational bottleneck to propagate the necessary information, **but one can observe computational bottlenecks even in short-range tasks**. An obvious example is an extremely high-degree node where, after *just one* layer, there is a high number of incoming messages, effectively creating a computational bottleneck in terms of information to be squashed into a fixed-size vector. Therefore, while the task of long-range is related to computational bottlenecks under classical DGNs, computational bottlenecks are not a prerogative of long-range tasks.

4.4 Belief: Topological Bottlenecks Associated with Long-range

The last belief we discuss is that topological bottlenecks are the primary obstacle to solving long-range tasks. This intuition stems from the fact that narrow cuts impede information flow between distant parts of the graph. It is indeed true that a topological bottleneck can worsen communication between distant nodes, especially if the bottleneck lies along a path that connects them. However, this perspective is limited in two important ways. First, a topological bottleneck is only harmful if it lies on the information paths between nodes that are supposed to communicate. A topological bottleneck may exist without affecting task-relevant dependencies. Second, the graph topology can worsen the long-range communication even in the absence of any identifiable topological bottleneck, by inducing computational bottlenecks. As we discussed in Section 4.1, the grid graph is an illustrative case: despite the lack of topological bottlenecks, to connect the opposite corner nodes we need, at least, as many message-passing layers as the distance between them, thus leading to a huge computational bottleneck that will likely hamper the ability to process long-range dependencies. In addition, some of the techniques that reduce topological bottlenecks rely on introducing more edges or nodes into the graph, with the aim of reducing the distance between far-away nodes that should communicate. It is important to note that, although these approaches might be beneficial for the task at hand, they also worsen the computational bottleneck by adding more branches to the computational tree.

In conclusion, this highlights a deeper issue: long-range problems are not solely caused by topological bottlenecks, rather they can be understood as a

form of information attenuation caused by computational bottlenecks in the message-passing mechanism, which can potentially be exacerbated by topological bottlenecks. Thus, solving a topological bottleneck is neither necessary nor sufficient to solve all long-range problems.

Message of the Section

- i) “Oversquashing” is an ambiguous term that led to unclear research statements. Talking about **computational and topological bottlenecks**, instead, better defines the research scope of a paper, since **they are two fundamentally distinct problems**.
- ii) There can be computational bottlenecks but no topological ones, and vice-versa. Hence, each of these two bottlenecks, though intertwined, deserves a dedicated research effort. This also means **creating ad-hoc benchmarks for each type of bottleneck rather than relying solely on real-world tasks**, where it is not as easy to distinguish the combined effect of the two bottlenecks.
- iii) Computational bottlenecks can happen in short-range as well as long-range tasks.
- iv) Performance issues in long-range tasks are not solely caused by topological bottlenecks; computational bottlenecks also play a role.

5 Conclusions

This position paper posits that the fast pace of advances in the graph machine learning field has generated several commonly accepted beliefs and hypotheses, rooted in the notions of oversmoothing, “oversquashing”, long-range tasks, and heterophily, which are the cause of misunderstandings between researchers. Our contribution was to highlight such beliefs in plain sight, provide an explanation for their emergence, and demystify them when necessary with simple counterexamples. First, we argued that OSM may not be an actual problem and that node embeddings’ separability should be preferred when looking for the root causes of performance degradation. Then, we showed how talking about computational and topological bottlenecks resolves most, if not all, inconsistencies generated by the inflated use of the “oversquashing” term. Finally, we highlighted the role of the task in statements involving homophily, heterophily, and long-range dependencies. By providing much-needed clarifications around these aspects, we hope to foster further advancements in the graph machine learning field.

Acknowledgments

A.A. acknowledges support from Intel corporation, a nominal grant received at the ELLIS Unit Alicante Foundation from the Regional Government of Valencia in Spain (Convenio Singular signed with Generalitat Valenciana, Conselleria de Innovación, Industria, Comercio y Turismo, Dirección General de Innovación)

and a grant by the Banc Sabadell Foundation. In addition, this work is partially funded by the European Union EU - HE ELIAS – Grant Agreement 101120237. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Health and Digital Executive Agency (HaDEA).

Bibliography

- [1] Abboud, R., Dimitrov, R., Ceylan, I.I.: Shortest path networks for graph property prediction. In: The 1st Learning on Graphs Conference (LoG) (2022)
- [2] Abu-El-Haija, S., Perozzi, B., Kapoor, A., Alipourfard, N., Lerman, K., Harutyunyan, H., Ver Steeg, G., Galstyan, A.: Mixhop: Higher-order graph convolutional architectures via sparsified neighborhood mixing. In: Proceedings of the 36th International Conference on Machine Learning (ICML) (2019)
- [3] Akansha, S.: Over-squashing in graph neural networks: A comprehensive survey. arXiv preprint arXiv:2308.15568 (2023)
- [4] Alon, U., Yahav, E.: On the bottleneck of graph neural networks and its practical implications. In: Proceedings of the 9th International Conference on Learning Representations (ICLR) (2021)
- [5] Arnaiz-Rodriguez, A., Begga, A., Escolano, F., Oliver, N.M.: DiffWire: Inductive Graph Rewiring via the Lovász Bound. In: Proceedings of the 1st Learning on Graphs Conference (LoG) (2022)
- [6] Arnaiz-Rodriguez, A., Vellingker, A.: Graph Learning: Principles, Challenges, and Open Directions. In: 41st International Conference on Machine Learning (ICML) (2024), tutorial
- [7] Arroyo, Á., Gravina, A., Gutteridge, B., Barbero, F., Gallicchio, C., Dong, X., Bronstein, M., Vandergheynst, P.: On vanishing gradients, over-smoothing, and over-squashing in gnns: Bridging recurrent and graph learning. arXiv preprint arXiv:2502.10818 (2025)
- [8] Attali, H., Buscaldi, D., Pernelle, N.: Delaunay graph: Addressing over-squashing and over-smoothing using delaunay triangulation. In: Proceedings of the 41st International Conference on Machine Learning (ICML) (2024)
- [9] Attali, H., Buscaldi, D., Pernelle, N.: Rewiring techniques to mitigate oversquashing and oversmoothing in gnns: A survey. arXiv preprint arXiv:2411.17429 (2024)
- [10] Bacciu, D., Errica, F., Micheli, A., Podda, M.: A gentle introduction to deep learning for graphs. *Neural Networks* **129** (2020)
- [11] Balla, J.: Over-squashing in riemannian graph neural networks. In: Proceedings of the 2nd Learning on Graphs Conference (LoG) (2023)
- [12] Banerjee, P.K., Karhadkar, K., Wang, Y.G., Alon, U., Montúfar, G.: Over-squashing in gnns through the lens of information contraction and graph expansion. In: Proceedings of the 58th Annual Allerton Conference on Communication, Control, and Computing (Allerton) (2022)
- [13] Barbero, F., Vellingker, A., Saberi, A., Bronstein, M.M., Giovanni, F.D.: Locality-aware graph rewiring in GNNs. In: Proceedings of the 12th International Conference on Learning Representations (ICLR) (2024)
- [14] Beaini, D., Huang, S., Cunha, J.A., Li, Z., Moisesescu-Pareja, G., Dymov, O., Maddrell-Mander, S., McLean, C., Wenkel, F., Müller, L., Mohamud, J.H.,

- Parviz, A., Craig, M., Koziarski, M., Lu, J., Zhu, Z., Gabellini, C., Klaser, K., Dean, J., Wognum, C., Sypetkowski, M., Rabusseau, G., Rabbany, R., Tang, J., Morris, C., Ravanelli, M., Wolf, G., Tossou, P., Mary, H., Bois, T., Fitzgibbon, A.W., Banaszewski, B., Martin, C., Masters, D.: Towards foundational models for molecular learning on large-scale multi-task datasets. In: Proceedings of the 12th International Conference on Learning Representations (ICLR) (2024)
- [15] Bi, W., Du, L., Fu, Q., Wang, Y., Han, S., Zhang, D.: Make heterophilic graphs better fit gnn: A graph rewiring approach. *IEEE Transactions on Knowledge and Data Engineering* (2024)
- [16] Black, M., Wan, Z., Nayyeri, A., Wang, Y.: Understanding oversquashing in gnns through the lens of effective resistance. In: Proceedings of the 40th International Conference on Machine Learning (ICML) (2023)
- [17] Bodnar, C., Di Giovanni, F., Chamberlain, B., Lio, P., Bronstein, M.: Neural sheaf diffusion: A topological perspective on heterophily and oversmoothing in gnns. In: Proceedings of the 36th Conference on Advances in Neural Information Processing Systems (NeurIPS) (2022)
- [18] Cai, C., Hy, T.S., Yu, R., Wang, Y.: On the connection between mpnn and graph transformer. In: Proceedings of the 40th International Conference on Machine Learning (ICML). pp. 3408–3430 (2023)
- [19] Cai, C., Wang, Y.: A note on over-smoothing for graph neural networks. In: Graph Representation Learning and Beyond Workshop, 37th International Conference on Machine Learning (ICML) (2020)
- [20] Castellana, D., Errica, F.: Investigating the interplay between features and structures in graph learning. In: MLG Workshop at ECML PKDD (2023)
- [21] Chandra, A.K., Raghavan, P., Ruzzo, W.L., Smolensky, R.: The electrical resistance of a graph captures its commute and cover times. In: Proceedings of the twenty-first annual ACM symposium on Theory of computing. pp. 574–586 (1989)
- [22] Chen, D., Lin, Y., Li, W., Li, P., Zhou, J., Sun, X.: Measuring and relieving the over-smoothing problem for graph neural networks from the topological view. In: Proceedings of the 34th AAAI conference on artificial intelligence (AAAI) (2020)
- [23] Chen, J., Zhu, J., Song, L.: Stochastic training of graph convolutional networks with variance reduction. In: Proceedings of the 35th International Conference on Machine Learning (ICML) (2018)
- [24] Chen, M., Wei, Z., Huang, Z., Ding, B., Li, Y.: Simple and deep graph convolutional networks. In: Proceedings of the 37th International Conference on Machine Learning (ICML) (2020)
- [25] Chen, T., Zhou, K., Duan, K., Zheng, W., Wang, P., Hu, X., Wang, Z.: Bag of tricks for training deeper graph neural networks: A comprehensive benchmark study. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **45**(3) (2022)
- [26] Chung, F.R.: Spectral graph theory, vol. 92. American Mathematical Soc. (1997)

- [27] Cong, W., Ramezani, M., Mahdavi, M.: On provable benefits of depth in training graph convolutional networks. *Advances in Neural Information Processing Systems* **34**, 9936–9949 (2021)
- [28] Deac, A., Lackenby, M., Veličković, P.: Expander graph propagation. In: *Proceedings of the 1st Learning on Graphs Conference (LoG)* (2022)
- [29] Di Giovanni, F., Giusti, L., Barbero, F., Luise, G., Lio, P., Bronstein, M.M.: On over-squashing in message passing neural networks: The impact of width, depth, and topology. In: *Proceedings of the 40th International Conference on Machine Learning (ICML)* (2023)
- [30] ud din, A.M., Qureshi, S.: Limits of depth: Over-smoothing and over-squashing in gnns. *Big Data Mining and Analytics* **7** (2024)
- [31] Dwivedi, V.P., Rampásek, L., Galkin, M., Parviz, A., Wolf, G., Luu, A.T., Beaini, D.: Long range graph benchmark. In: *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS)* (2022)
- [32] Ellens, W., Spieksma, F.M., Van Mieghem, P., Jamakovic, A., Kooij, R.E.: Effective graph resistance. *Linear algebra and its applications* **435**(10), 2491–2506 (2011)
- [33] Epping, B., René, A., Helias, M., Schaub, M.T.: Graph neural networks do not always oversmooth. In: *Proceedings of the 38th Conference on Neural Information Processing Systems (NeurIPS)* (2024)
- [34] Errica, F.: On class distributions induced by nearest neighbor graphs for node classification of tabular data. In: *Proceedings of the 37th Conference on Advances in Neural Information Processing Systems (NeurIPS)* (2023)
- [35] Errica, F., Christiansen, H., Zaverkin, V., Maruyama, T., Niepert, M., Alesiani, F.: Adaptive message passing: A general framework to mitigate oversmoothing, oversquashing, and underreaching. In: *Proceedings of the 42nd International Conference on Machine Learning (ICML)* (2025)
- [36] Fesser, L., Weber, M.: Mitigating over-smoothing and over-squashing using augmentations of forman-ricci curvature. In: *Proceedings of the 4th Learning on Graphs Conference (LoG)* (2024)
- [37] Gabrielsson, R.B., Yurochkin, M., Solomon, J.: Rewiring with positional encodings for graph neural networks. *Transactions on Machine Learning Research* (2023)
- [38] Gasteiger, J., Bojchevski, A., Günnemann, S.: Predict then propagate: Graph neural networks meet personalized pagerank. In: *Proceedings of the 7th International Conference on Learning Representations (ICLR)* (2019)
- [39] Gilmer, J., Schoenholz, S.S., Riley, P.F., Vinyals, O., Dahl, G.E.: Neural message passing for quantum chemistry. In: *Proceedings of the 34th International Conference on Machine Learning (ICML)* (2017)
- [40] Giovanni, F.D., Rowbottom, J., Chamberlain, B.P., Markovich, T., Bronstein, M.M.: Understanding convolution on graphs via energies. *Transactions on Machine Learning Research* (2023)
- [41] Giovanni, F.D., Rusch, T.K., Bronstein, M., Deac, A., Lackenby, M., Mishra, S., Veličković, P.: How does over-squashing affect the power of GNNs? *Transactions on Machine Learning Research* (2024)

- [42] Giraldo, J.H., Skianis, K., Bouwmans, T., Malliaros, F.D.: On the trade-off between over-smoothing and over-squashing in deep graph neural networks. In: Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (ICKM) (2023)
- [43] Gong, C., Cheng, Y., Yu, J., Xu, C., Shan, C., Luo, S., Li, X.: A survey on learning from graphs with heterophily: Recent advances and future directions. arXiv preprint arXiv:2401.09769 (2024)
- [44] Gori, M., Monfardini, G., Scarselli, F.: A new model for learning in graph domains. In: Proceedings of the International Joint Conference on Neural Networks (IJCNN) (2005)
- [45] Gravina, A., Eliasof, M., Gallicchio, C., Bacciu, D., Schönlieb, C.B.: On oversquashing in graph neural networks through the lens of dynamical systems. In: Proceedings of the 39th AAAI Conference on Artificial Intelligence (AAAI) (2025)
- [46] Gutteridge, B., Dong, X., Bronstein, M.M., Di Giovanni, F.: Drew: Dynamically rewired message passing with delay. In: Proceedings of the 40th International Conference on Machine Learning (ICML) (2023)
- [47] Hamilton, W.L.: Graph representation learning. Morgan & Claypool Publishers (2020)
- [48] Hasanzadeh, A., Hajiramezanali, E., Boluki, S., Zhou, M., Duffield, N., Narayanan, K., Qian, X.: Bayesian graph neural networks with adaptive connection sampling. In: Proceedings of the 37th International Conference on Machine Learning (ICML) (2020)
- [49] Hoang, N., Maehara, T., Murata, T.: Revisiting graph neural networks: Graph filtering perspective. In: Proceedings of the 25th International Conference on Pattern Recognition (ICPR) (2021)
- [50] Huang, K., Wang, Y.G., Li, M., Lio, P.: How universal polynomial bases enhance spectral graph neural networks: Heterophily, over-smoothing, and over-squashing. In: Proceedings of the 41st International Conference on Machine Learning (ICML) (2024)
- [51] Huang, W., Rong, Y., Xu, T., Sun, F., Huang, J.: Tackling over-smoothing for general graph convolutional networks. arXiv preprint arXiv:2008.09864 (2020)
- [52] Hwang, E., Thost, V., Dasgupta, S.S., Ma, T.: An analysis of virtual nodes in graph neural networks for link prediction (extended abstract). In: Proceedings of the 1st Learning on Graphs Conference (LoG) (2022)
- [53] Jamadandi, A., Rubio-Madrigal, C., Burkholz, R.: Spectral graph pruning against over-squashing and over-smoothing. In: Proceedings of the 38th Conference on Neural Information Processing Systems (NeurIPS) (2024)
- [54] Karhadkar, K., Banerjee, P.K., Montufar, G.: FoSR: First-order spectral rewiring for addressing oversquashing in GNNs. In: Proceedings of the 11th International Conference on Learning Representations (ICLR) (2023)
- [55] Keriven, N.: Not too little, not too much: a theoretical analysis of graph (over) smoothing. In: Proceedings of the 36th Conference on Advances in Neural Information Processing Systems (NeurIPS) (2022)

- [56] Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: Proceedings of the 5th International Conference on Learning Representations (ICLR) (2017)
- [57] Klein, D.J., Randić, M.: Resistance distance. *Journal of mathematical chemistry* **12**, 81–95 (1993)
- [58] Li, G., Muller, M., Thabet, A., Ghanem, B.: Deepgcns: Can gcns go as deep as cnns? In: Proceedings of the 17th IEEE/CVF International Conference on Computer Vision (ICCV) (2019)
- [59] Li, Q., Han, Z., Wu, X.M.: Deeper insights into graph convolutional networks for semi-supervised learning. In: Proceedings of the 32nd AAAI conference on artificial intelligence (AAAI) (2018)
- [60] Lim, D., Hohne, F., Li, X., Huang, S.L., Gupta, V., Bhalerao, O., Lim, S.N.: Large scale learning on non-homophilous graphs: New benchmarks and strong simple methods. In: Proceedings of the 35th Conference on Neural Information Processing Systems (NeurIPS) (2021)
- [61] Liu, M., Gao, H., Ji, S.: Towards deeper graph neural networks. In: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD). pp. 338–348 (2020)
- [62] Liu, Y., Zhou, C., Pan, S., Wu, J., Li, Z., Chen, H., Zhang, P.: Curvdrop: A ricci curvature based approach to prevent graph neural networks from over-smoothing and over-squashing. In: Proceedings of the ACM International World Wide Web Conference (WWW) (2023)
- [63] Liu, Y., Zheng, Y., Zhang, D., Lee, V.C., Pan, S.: Beyond smoothing: Unsupervised graph representation learning with edge heterophily discriminating. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI) (2023)
- [64] Longa, A., Lachi, V., Santin, G., Bianchini, M., Lepri, B., Lio, P., Francoscelli, Passerini, A.: Graph neural networks for temporal graphs: State of the art, open challenges, and opportunities. *Transactions on Machine Learning Research* (2023)
- [65] Lovász, L.: Random walks on graphs. *Combinatorics*, Paul erdos is eighty **2**(1-46), 4 (1993), <https://web.cs.elte.hu/~lovasz/erdos.pdf>
- [66] Luan, S., Hua, C., Lu, Q., Zhu, J., Zhao, M., Zhang, S., Chang, X.W., Precup, D.: Is heterophily a real nightmare for graph neural networks to do node classification? arXiv preprint arXiv:2109.05641 (2021)
- [67] Luan, S., Hua, C., Lu, Q., Zhu, J., Zhao, M., Zhang, S., Chang, X.W., Precup, D.: Revisiting heterophily for graph neural networks. In: Proceedings of the 36th Conference on Advances in Neural Information Processing Systems (NeurIPS) (2022)
- [68] Luan, S., Hua, C., Xu, M., Lu, Q., Zhu, J., Chang, X.W., Fu, J., Leskovec, J., Precup, D.: When do graph neural networks help with node classification? investigating the homophily principle on node distinguishability. In: Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS) (2023)
- [69] Ma, Y., Liu, X., Shah, N., Tang, J.: Is homophily a necessity for graph neural networks? In: 10th International Conference on Learning Representations (ICLR) (2022)

- [70] Maskey, S., Paolino, R., Bacho, A., Kutyniok, G.: A fractional graph laplacian approach to oversmoothing. In: Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS) (2023)
- [71] McPherson, M., Smith-Lovin, L., Cook, J.M.: Birds of a feather: Homophily in social networks. *Annual review of sociology* **27** (2001)
- [72] Micheli, A.: Neural network for graphs: A contextual constructive approach. *IEEE Transactions on Neural Networks* **20** (2009)
- [73] Micheli, A., Sestito, A.: A new neural network model for contextual processing of graphs. In: Proceedings of the Italian Workshop on Neural Networks (WIRN) (2005)
- [74] Morris, C., Lipman, Y., Maron, H., Rieck, B., Kriege, N.M., Grohe, M., Fey, M., Borgwardt, K.: Weisfeiler and leman go machine learning: The story so far. *Journal of Machine Learning Research* **24** (2023)
- [75] Namata, G.M., London, B., Getoor, L., Huang, B.: Query-driven active surveying for collective classification. In: Proceedings of the Workshop on Mining and Learning with Graphs (MLG) (2012)
- [76] Nguyen, K., Hieu, N.M., Nguyen, V.D., Ho, N., Osher, S., Nguyen, T.M.: Revisiting over-smoothing and over-squashing using ollivier-ricci curvature. In: Proceedings of the 40th International Conference on Machine Learning (ICML) (2023)
- [77] Oono, K., Suzuki, T.: Graph neural networks exponentially lose expressive power for node classification. In: Proceedings of the 8th International Conference on Learning Representations (ICLR) (2020)
- [78] Park, M., Choi, S., Heo, J., Park, E., Kim, D.: The oversmoothing fallacy: A misguided narrative in gnn research. *arXiv preprint arXiv:2506.04653* (2025)
- [79] Pei, H., Wei, B., Chang, K.C.C., Lei, Y., Yang, B.: Geom-gcn: Geometric graph convolutional networks. In: Proceedings of the 8th International Conference on Learning Representations (ICLR) (2020)
- [80] Platonov, O., Kuznedelev, D., Babenko, A., Prokhorenkova, L.: Characterizing graph datasets for node classification: Homophily-heterophily dichotomy and beyond. In: Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS) (2023)
- [81] Platonov, O., Kuznedelev, D., Diskin, M., Babenko, A., Prokhorenkova, L.: A critical look at the evaluation of GNNs under heterophily: Are we really making progress? In: Proceedings of the 1th International Conference on Learning Representations (ICLR) (2023)
- [82] Qiu, H., Hancock, E.R.: Clustering and embedding using commute times. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **29**(11), 1873–1890 (2007)
- [83] Rong, Y., Huang, W., Xu, T., Huang, J.: Dropedge: Towards deep graph convolutional networks on node classification. In: Proceedings of the 8th International Conference on Learning Representations (2020)
- [84] Roth, A., Liebig, T.: Rank collapse causes over-smoothing and over-correlation in graph neural networks. In: Proceedings of the 3rd Learning on Graphs Conference (LoG) (2024)

- [85] Rusch, T.K., Bronstein, M.M., Mishra, S.: A survey on oversmoothing in graph neural networks. *arXiv preprint arXiv:2303.10993* (2023)
- [86] Rusch, T.K., Chamberlain, B., Rowbottom, J., Mishra, S., Bronstein, M.: Graph-coupled oscillator networks. In: *Proceedings of the 39th International Conference on Machine Learning (ICML)* (2022)
- [87] Scarselli, F., Gori, M., Tsoi, A.C., Hagenbuchner, M., Monfardini, G.: The graph neural network model. *IEEE Transactions on Neural Networks* **20** (2009)
- [88] Sen, P., Namata, G., Bilgic, M., Getoor, L., Galligher, B., Eliassi-Rad, T.: Collective classification in network data. *AI magazine* **29**(3), 93–93 (2008)
- [89] Shao, Z., Shi, D., Han, A., Guo, Y., Zhao, Q., Gao, J.: Unifying oversmoothing and over-squashing in graph neural networks: A physics informed approach and beyond. *arXiv preprint arXiv:2309.02769* (2023)
- [90] Shi, D., Han, A., Lin, L., Guo, Y., Gao, J.: Exposition on over-squashing problem on gnns: Current methods, benchmarks and challenges. *arXiv preprint arXiv:2311.07073* (2023)
- [91] Southern, J., Di Giovanni, F., Bronstein, M., Lutzeyer, J.F.: Understanding virtual nodes: Oversquashing and node heterogeneity. In: *International Conference on Learning Representations (ICLR)* (2025)
- [92] Sperduti, A., Starita, A.: Supervised neural networks for the classification of structures. *IEEE Transactions on Neural Networks* **8**(3) (1997)
- [93] Stanovic, S., Gaüzère, B., Brun, L.: Graph neural networks with maximal independent set-based pooling: Mitigating over-smoothing and over-squashing. *Pattern Recognition Letters* **187** (2025)
- [94] Sun, Q., Li, J., Yuan, H., Fu, X., Peng, H., Ji, C., Li, Q., Yu, P.S.: Position-aware structure learning for graph topology-imbalance by relieving under-reaching and over-squashing. In: *Proceedings of the 31st ACM International Conference on Information & Knowledge Management (ICKM)*. pp. 1848–1857 (2022)
- [95] Topping, J., Giovanni, F.D., Chamberlain, B.P., Dong, X., Bronstein, M.M.: Understanding over-squashing and bottlenecks on graphs via curvature. In: *Proceedings of the 10th International Conference on Learning Representations (ICLR)* (2022)
- [96] Tori, F., Holst, V., Ginis, V.: The effectiveness of curvature-based rewiring and the role of hyperparameters in GNNs revisited. In: *Proceedings of the 13th International Conference on Learning Representations (ICLR)* (2025)
- [97] Tortorella, D., Micheli, A.: Leave graphs alone: Addressing over-squashing without rewiring. In: *The 1st Learning on Graphs Conference (LoG)* (2022)
- [98] Wang, H., Leskovec, J.: Combining graph convolutional neural networks and label propagation. *ACM Transactions on Information Systems* **40**(4) (2021)
- [99] Wang, J., Guo, Y., Yang, L., Wang, Y.: Understanding heterophily for graph neural networks. In: *Proceedings of the 41st International Conference on Machine Learning (ICML)* (2024)
- [100] Wang, K., Yang, Y., Saha, I., Allen-Blanchette, C.: Understanding oversmoothing in gnns as consensus in opinion dynamics. *arXiv preprint arXiv:2501.19089* (2025)

- [101] Wu, X., Ajorlou, A., Wu, Z., Jadbabaie, A.: Demystifying oversmoothing in attention-based graph neural networks. In: Proceedings of the 37th Conference on Advances in Neural Information Processing Systems (NeurIPS) (2023)
- [102] Wu, X., Chen, Z., Wang, W.W., Jadbabaie, A.: A non-asymptotic analysis of oversmoothing in graph neural networks. In: Proceedings of the 11th International Conference on Learning Representations (ICLR) (2023)
- [103] Xu, K., Hu, W., Leskovec, J., Jegelka, S.: How powerful are graph neural networks? In: Proceedings of the 7th International Conference on Learning Representations (ICLR) (2019)
- [104] Xu, K., Li, C., Tian, Y., Sonobe, T., Kawarabayashi, K.i., Jegelka, S.: Representation learning on graphs with jumping knowledge networks. In: Proceedings of the 35th International Conference on Machine Learning (ICML) (2018)
- [105] Yang, C., Wang, R., Yao, S., Liu, S., Abdelzaher, T.: Revisiting oversmoothing in deep gcns. arXiv preprint arXiv:2003.13663 (2020)
- [106] Yu, W., Ma, X., Bailey, J., Zhan, Y., Wu, J., Du, B., Hu, W.: Graph structure reforming framework enhanced by commute time distance for graph classification. *Neural Networks* **168** (2023)
- [107] Zhang, K., Deidda, P., Higham, D., Tudisco, F.: Rethinking oversmoothing in graph neural networks: A rank-based perspective. arXiv preprint arXiv:2502.04591 (2025)
- [108] Zhang, X., Wang, L., Helwig, J., Luo, Y., Fu, C., Xie, Y., Liu, M., Lin, Y., Xu, Z., Yan, K., et al.: Artificial intelligence for science in quantum, atomistic, and continuum systems. arXiv preprint arXiv:2307.08423 (2023)
- [109] Zhao, L., Akoglu, L.: Pairnorm: Tackling oversmoothing in gnns. In: Proceedings of the 8th International Conference on Learning Representations (ICLR) (2020)
- [110] Zheng, Y., Luan, S., Chen, L.: What is missing in homophily? disentangling graph homophily for graph neural networks. arXiv preprint arXiv:2406.18854 (2024)
- [111] Zhou, K., Huang, X., Li, Y., Zha, D., Chen, R., Hu, X.: Towards deeper graph neural networks with differentiable group normalization. In: Proceedings of the 34th Conference on Advances in Neural Information Processing Systems (NeurIPS) (2020)
- [112] Zhou, K., Huang, X., Zha, D., Chen, R., Li, L., Choi, S.H., Hu, X.: Dirichlet energy constrained learning for deep graph neural networks. In: Proceedings of the 35th Conference on Advances in Neural Information Processing Systems (NeurIPS) (2021)
- [113] Zhu, J., Yan, Y., Zhao, L., Heimann, M., Akoglu, L., Koutra, D.: Beyond homophily in graph neural networks: Current limitations and effective designs. In: Proceedings of the 34th Conference on Neural Information Processing Systems (NeurIPS) (2020)

A Background

A.1 Deep Graph Networks

We provide a brief excursus into Deep Graph Networks for readers new to the topic.

We can define a graph as a tuple $g = (\mathcal{V}_g, \mathcal{E}_g, \mathcal{X}_g, \mathcal{A}_g)$, with \mathcal{V}_g the set of nodes, \mathcal{E}_g the set of edges (oriented or not oriented) connecting pairs of nodes. \mathcal{E}_g encodes the topological information of the graph and can be represented as an adjacency matrix: a binary square matrix \mathbf{A} where \mathbf{A}_{uv} is 1 if there is an edge between u and v , and it is 0 otherwise. Additional node and edge features belong are represented by $\mathbf{x}_v \in \mathcal{X}_g$ and $\mathbf{a}_{uv} \in \mathcal{A}_g$, respectively. \mathcal{X}_g can be, for instance, $\mathbb{R}^d, d \in \mathbb{N}^+$.

The neighborhood of a node v is the set of nodes that are connected to v by an oriented edge, i.e., $\mathcal{N}_v = \{u \in \mathcal{V}_g | (u, v) \in \mathcal{E}_g\}$. If the graph is undirected, we convert each non-oriented edge into two oriented but opposite ones.

The main mechanism of DGNs is the repeated aggregation of neighbors' information, which gives rise to the spreading of local information across the graph. The process is simple: i) at iteration ℓ , each node receives "messages" (usually just node representations) from the neighbors and processes them into a single new message; ii) the message is used to update the representation of that node. Both steps involve learnable functions, so DGNs can learn to capture the relevant correlations in the graph.

Most DGNs implement a synchronous message-passing mechanism, meaning each node always receives information from all neighbors at every iteration step. This local and iterative processing is at the core of DGNs' efficiency since computation can be easily parallelized across nodes. In addition, being local means being independent of the graph's size. When one learns the same function for all message passing iterations, we talk about *recurrent* architectures, as the GNN of [44, 87]; on the contrary, when one learns a separate parametrization for a finite number of iterations (also known as layers), we talk about *convolutional* architectures as the NN4G of [72, 73].

The neighborhood aggregation is usually implemented using permutation-invariant functions, which make learning possible on cyclic graphs that have no consistent topological ordering of their nodes. A rather general and classical neighborhood aggregation mechanism for node v at layer/step $\ell+1$ is the following:

$$\mathbf{h}_v^{\ell+1} = \phi^{\ell+1} \left(\mathbf{h}_v^\ell, \Psi(\{\psi^{\ell+1}(\mathbf{h}_u^\ell) \mid u \in \mathcal{N}_v\}) \right) \quad (2)$$

where \mathbf{h}_u^ℓ is the node embedding of u at layer/step ℓ , ϕ and ψ implement learnable functions, and Ψ is a permutation invariant aggregation function. Note that $\mathbf{h}_v^0 = \mathbf{x}_v$. For instance, the Graph Convolutional Network of [56] implements the following aggregation, which is a special case of the above equation:

$$\mathbf{h}_v^{\ell+1} = \sigma(\mathbf{W}^{\ell+1} \sum_{u \in \mathcal{N}(v)} \hat{\mathbf{L}}_{uv} \mathbf{h}_u^\ell), \quad (3)$$

with $\hat{\mathbf{L}}$ being the normalized graph Laplacian, \mathbf{W} is a learnable weight matrix and σ is a non-linear activation function.

A.2 Some OSM Definitions

In order to keep the paper self-contained and to illustrate the different ways OSM has been defined, we briefly review the most common definitions. The following subsection summarizes the most commonly used OSM metrics. Some of these metrics are used in the main text.

Let $D = \text{diag}(d_1, \dots, d_n)$ be the degree matrix with $d_u = \sum_v A_{uv}$. The combinatorial graph Laplacian is defined as $\mathbf{L} = \mathbf{D} - \mathbf{A}$, with eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. In addition, the symmetric normalized Laplacian, defined by $\hat{\mathbf{L}} = \mathbf{I} - \mathbf{D}^{1/2} \mathbf{A} \mathbf{D}^{1/2}$, normalizes each entry to remove the influence of node degrees and its spectrum lies in $[0, 2]$.

DE [26] was used to measure OSM in [19] is defined as

$$\text{DE}(\mathbf{H}^\ell) = \text{Tr}((\mathbf{H}^\ell)^T \hat{\mathbf{L}} \mathbf{H}^\ell) = \frac{1}{2} \sum_{u,v \in \mathcal{E}} \left\| \frac{\mathbf{h}_u}{\sqrt{d_u}} - \frac{\mathbf{h}_v}{\sqrt{d_v}} \right\|_2^2 \quad (4)$$

Note that DE can also be computed using \mathbf{L}

$$\text{DE}(\mathbf{H}^\ell) = \text{Tr}((\mathbf{H}^\ell)^T \mathbf{L} \mathbf{H}^\ell) = \frac{1}{2} \sum_{u,v \in \mathcal{E}} \|\mathbf{h}_u - \mathbf{h}_v\|_2^2 \quad (5)$$

Other OSM metrics include with respect to the norm of node embeddings. For instance, Rayleigh Quotient (RQ) [26], which can be seen as normalized DE, is defined by

$$\text{RQ} = \frac{\text{Tr}((\mathbf{H}^\ell)^T \hat{\mathbf{L}} \mathbf{H}^\ell)}{\|\mathbf{H}^\ell\|_F^2} \quad (6)$$

first proposed for OSM analysis in [19] and later used by [40, 70, 84, 105]. RQ discerns whether embeddings become *relatively* smoother, independent of magnitude.

Mean Absolute Deviation (MAD) [22] averages cosine dissimilarity between a node and its neighbors.

$$\text{MAD}_G(\mathbf{H}^\ell) = \frac{1}{n} \sum_{v \in \mathcal{V}} \sum_{u \in \mathcal{N}_v} 1 - \frac{\mathbf{h}_v^{\ell T} \mathbf{h}_u^\ell}{\|\mathbf{h}_v^\ell\| \|\mathbf{h}_u^\ell\|} \quad (7)$$

The smoothness metric SMV [61] captures a global node-distance average over all node pairs:

$$\text{SMV} = \frac{1}{n} \sum_{u \in \mathcal{V}} \frac{1}{n-1} \sum_{v \neq u \in \mathcal{V}} \frac{1}{2} \left\| \frac{\mathbf{h}_u}{\|\mathbf{h}_u\|} - \frac{\mathbf{h}_v}{\|\mathbf{h}_v\|} \right\| \quad (8)$$

In the main text, we primarily consider DE and RQ. However, all notions convey the same intuition, loss of discriminative variation, yet, as we show, can disagree in practice.

Convergence-rate results. Known theoretical bounds show $\text{DE}(H^k)$ decays exponentially with depth k (e.g., the rate depends on weight spectra and graph eigenvalues) [51, 77]. Such convergence rate results primarily make assumptions on the architecture, weight matrix, and activation functions, and may be altered by skip connections, normalization layers, or simple rescaling such as $2W$ [84].

For instance [19] propose a bound on the DE of two consecutive message passing layers (similar bounds found in [77, 112])

$$\hat{\text{DE}}(\mathbf{H}^\ell) \leq (1 - \lambda_2)^2 s_{\max}^\ell \hat{\text{DE}}(\mathbf{H}^{\ell-1})$$

being s_{\max}^ℓ the square of the maximum singular value of W^ℓ , and λ_2 the second smallest eigenvalue of the Laplacian, i.e. the spectral gap. The proof holds when $s_{\max}^\ell < 1/(1 - \lambda_2)$.

In addition, [40] further relate Laplacian eigenvalues and weight spectra to predict whether RQ converges to 0 (collapse) or to λ_{\max} (no collapse).

A.3 Some OSQ Definitions

For completeness, we summarize some of the most commonly used metrics in the OSQ literature and their relationships. These quantities mainly measure three different aspects of the graph: (i) *sensitivity/Jacobian* measures that capture how information from a distant node u affects a target node v after K message-passing layers; (ii) *topological bottleneck* proxies such as Cheeger-type cut ratios, graph spectrum or curvature scores; and (iii) *distance-based* quantities such as effective resistance that upper-bound information flow.

Sensitivity For a K -layer GNN let $\mathbf{h}_u^{(k)}$ denote the embedding of node u at layer k . A first proxy for OSQ is the *Influence Score* of [104],

$$I(u, v) = \left\| \frac{\partial \mathbf{h}_u^K}{\partial \mathbf{h}_v^0} \right\|. \quad (9)$$

[16] sum these sensitivities over *all* unordered pairs,

$$\sum_{u \neq v \in V} \left\| \frac{\partial \mathbf{h}_u^K}{\partial \mathbf{h}_v^0} \right\|, \quad (10)$$

to obtain a graph-level indicator of how much information is lost.

[40] propose a *symmetric Jacobian obstruction* that removes self-influence and degree bias. They define the Jacobian obstruction of node v with respect to node u at layer m as

$$\tilde{\mathbf{J}}_k^{(m)}(v, u) := \left(\frac{1}{d_v} \frac{\partial \mathbf{h}_v^{(m)}}{\partial \mathbf{h}_v^{(k)}} - \frac{1}{\sqrt{d_v d_u}} \frac{\partial \mathbf{h}_v^{(m)}}{\partial \mathbf{h}_u^{(k)}} \right) + \left(\frac{1}{d_u} \frac{\partial \mathbf{h}_u^{(m)}}{\partial \mathbf{h}_u^{(k)}} - \frac{1}{\sqrt{d_v d_u}} \frac{\partial \mathbf{h}_u^{(m)}}{\partial \mathbf{h}_v^{(k)}} \right), \quad (11)$$

being the extension to the Jacobian obstruction of node v with respect to node u after m layers defined as

$$\tilde{O}^m(u, v) = \sum_{k=0}^m \left\| \tilde{\mathbf{J}}_k^{(m)}(v, u) \right\|. \quad (12)$$

Topological Bottlenecks Many OSQ papers measure topological (structural) bottlenecks using spectral or curvature quantities. Note that here we give an intuition based on the spectral metrics [5, 12, 54], but a significant part of the literature uses metrics based on curvature [42, 62, 76, 95].

First, the topological bottleneck can be measured by Cheeger Constant [26], which is the size of the min-cut of the graph.

$$h_G = \min_{S \subset V} \frac{|\{e = (u, v) : u \in S, v \in \bar{S}\}|}{\min\{\text{vol}(S), \text{vol}(\bar{S})\}}$$

A small h_G means one can separate G into two large-volume parts by removing only a few edges, i.e. a severe *topological bottleneck*. Cheeger's inequality links h_G to the spectrum of G :

$$\frac{h_G^2}{2} \leq \lambda_2 \leq 2h_G,$$

where λ_2 is the second eigenvalue of the normalized Laplacian.

Pairwise Distances The commute time between two nodes [65] is defined as the expected number of steps that a random walker needs to go from node u to v and come back. The Effective Resistance between two nodes [21], R_{uv} , is the commute time divided by the volume of the graph [57], which is the sum of the degrees of all nodes in the graph. The effective resistance between two nodes is computed as

$$R_{u,v} = L_{ii}^+ + L_{jj}^+ - 2L_{ij}^+$$

being $\mathbf{L}^+ = \sum_{i>0} \frac{1}{\lambda_i} \phi_i \phi_i^T$ the pseudoinverse of \mathbf{L}

Then, some measures derived from this metric can be connected with the topological bottleneck [21, 26, 82]. For instance, the maximum effective resistance of a graph is connected with the Cheeger constant as per $R_{\max} = \max_{u,v \in V} R_{uv}$

$$R_{\max} \leq \frac{1}{h_G^2}$$

and thus also bounded by the spectral gap as

$$\frac{1}{n\lambda_2} \leq R_{\max} \leq \frac{2}{\lambda_2}.$$

In addition, the total effective resistance $R_{\text{tot}} = 1/2 \sum_{u,v \in V} R_{uv}$ is bounded to the spectral gap [32]:

$$\frac{n}{\lambda_2} \leq R_{\text{tot}} \leq \frac{n(n-1)}{\lambda_2}$$

Note that the total effective resistance also equals the sum of the spectrum of \mathbf{L}^+ $R_{\text{tot}} = n \sum_2^n 1/\lambda_n$.

Connecting Sensitivity and Topological Distances The larger Total Effective Resistance ($R_{tot} = \sum_{(u,v) \in V} R_{uv}$) is, the lower the sum of pairwise sensitivities [16]:

$$\sum_{u,v \in V \times V} \left\| \frac{\partial h_v^{(r)}}{\partial h_u^{(0)}} \right\| \leq c(b - R_{tot}) \quad (13)$$

The larger the Effective Resistance is, the higher the Symmetric Jacobian Obstruction [40]:

$$\tilde{O}^m(u, v) = \sum_{k=0}^m \left\| \tilde{\mathbf{J}}_k^{(m)}(v, u) \right\| \leq c R_{u,v} \quad (14)$$

A.4 Computational Bottleneck

Oversquashing can also be seen through the perspective of the message-passing computational graph: each message-passing layer expands the set of nodes whose features can influence a target node. If this *receptive field* grows fast, any fixed-width DGN “squash” many signals into a single vector.

Receptive Field Following [4, 23] the receptive field was defined recursively as:

$$\mathcal{N}_v^K := \mathcal{N}_v^{K-1} \cup \{w \mid w \in \mathcal{N}_u \wedge u \in \mathcal{N}_v^{K-1}\} \quad \text{and} \quad \mathcal{N}_v^1 = \mathcal{N}_v \quad (15)$$

which can be also seen as the set of K -hop neighbors neighbors, i.e. nodes that are reachable from v within K hops. The number of nodes in each node’s receptive field can grow exponentially with the number of layers $|\mathcal{N}_v^K| = \mathcal{O}(\exp(K))$ [23]. For instance, in a rooted binary tree each layer has exactly b^{K-1} new neighbors, so $|\mathcal{N}_v^K| = 1 + b + b^2 + \dots + b^{K-1} = \Theta(b^K)$.

When evaluating the actual *computational* graph resulting from message-passing, duplicates matter: a node that appears in several branches of the computation tree contributes multiple times, since each distinct walk contributes a separate message. We therefore use the multiset notation to define the computational tree for a node v :

$$\mathcal{M}_v^1 := \mathcal{N}_v, \quad \mathcal{M}_v^K := \mathcal{M}_v^{K-1} \uplus \left\{ \biguplus_{u \in \mathcal{M}_v^{K-1}} \mathcal{N}_u \right\}. \quad (16)$$

Therefore, we can define the notion of an "exponentially-growing receptive field" [4] as follows.

Definition 1 (Computational Bottleneck). *For a given node v and number of message passing layers K , the computational bottleneck of node v is defined as $|\mathcal{M}_v^K|$.*

The size of the computational bottleneck (multiset receptive field) at node v , can be computed as:

$$|\mathcal{M}_v^K| := \sum_{\ell=1}^K \|A^\ell[v, :]\|_1 = \sum_{\ell=1}^K \sum_{u \in \mathcal{V}} (A^\ell)_{u,v} \quad (17)$$

This definition counts every distinct length- ℓ walk from v to any node u . Equation (17) is exactly the row-sum of the powers of the adjacency matrix; it therefore matches the size of the *computational tree*.

Note that the size of the set-based receptive field corresponds to the support of the multiset \mathcal{M}_v^K , denoted $\mathcal{N}_v^K := \text{supp}(\mathcal{M}_v^K)$. Therefore, the multiset size $|\mathcal{M}_v^K|$ is always greater than or equal to the size of the support, $|\mathcal{M}_v^K| \geq |\mathcal{N}_v^K|$, since it accounts for path multiplicity.

In early deep-graph networks literature, [72] introduced the idea by using the term “contextual window”: deeper layers aggregate exponentially many paths unless skip connections or global pooling curb the growth. The multiset perspective in Eq. (17) makes this explosion explicit and the matrix computation directly links to matrix-power interpretations of message passing.

Figure 5 visualises the difference between the set size $|\mathcal{N}_v^K|$ and the multiset size $|\mathcal{M}_v^K|$ on a toy graph and on a stochastic block model (SBM).

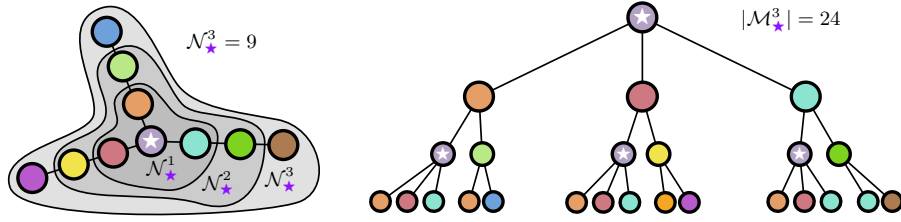


Fig. 5: **Computational Bottleneck.** Illustration of a receptive field – defined with sets (K -hop neighborhood) – and the definition of computational bottleneck measured as the size of the computational graph – defined with multisets.

In conclusion, we note that in message-passing, the computational bottleneck is driven not by how many distinct vertices are in the K -hop neighborhood, but by the size of the computational graph.

B Sensitivity Decreases on a Grid Graph without Topological Bottlenecks

To show that low sensitivity does not necessarily imply a topological bottleneck, Figure 6 analyzes sensitivity’s decreasing trend when the number of message passing layers L increases on the grid graph of Figure 4 of size 10×10 . Increasing the size of the embedding space postpones the collapse of the sensibility.

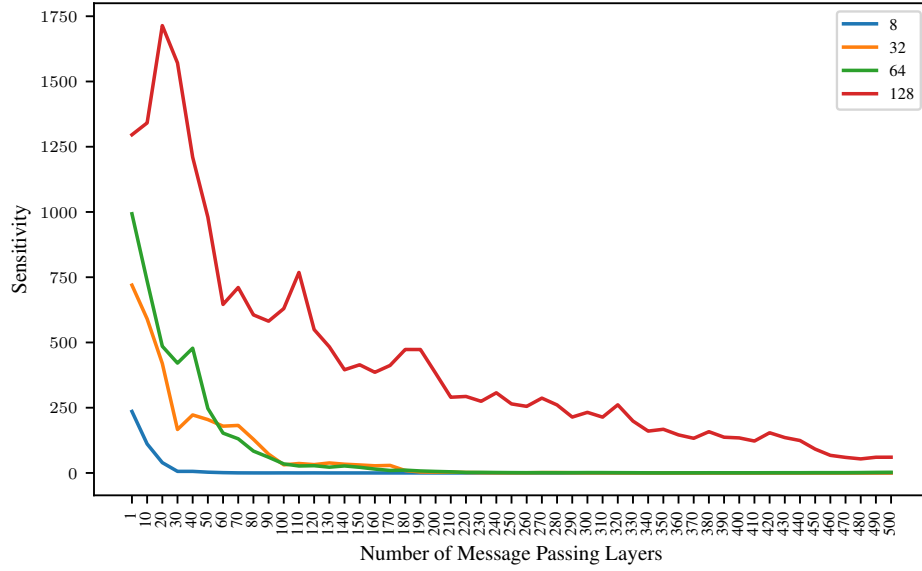


Fig.6: We plot the sensitivity of the grid graph of Figure 4 for the Graph Convolutional Network [56] model for different node embedding sizes.

C Oversmoothing Does Not Always Happen

For enhanced clarity and to allow for a more detailed examination of the OSM behavior discussed in Section 2, a larger version of Figure 1 is provided in Figure 7.

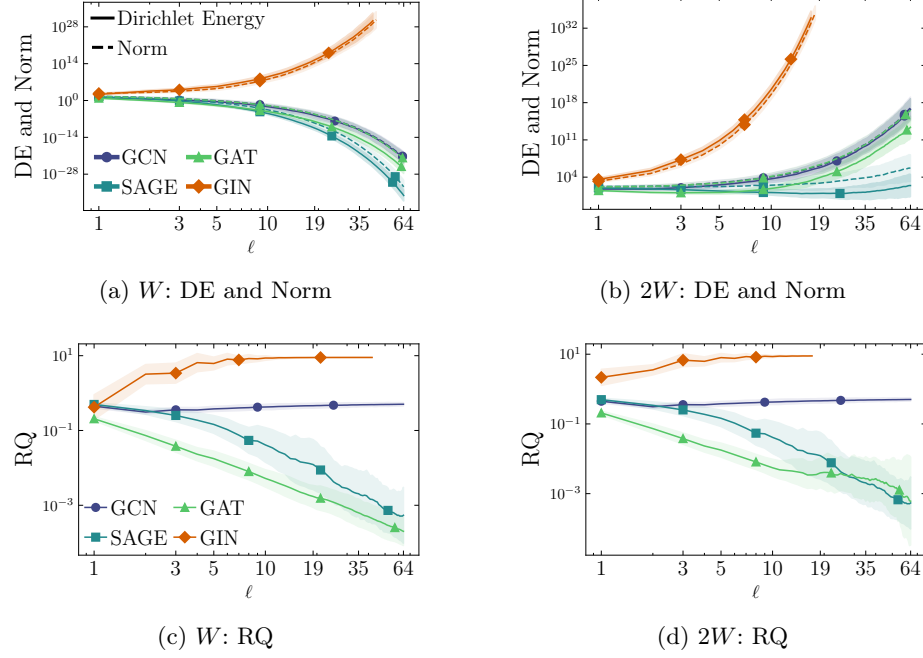


Fig. 7: **Larger version of Figure 1.** (a-b): We depict the evolution, with increasing number of layers, of the $DE = \text{tr}(X^T \Delta X)$ and the feature norm $\|X\|_F$, using W and $2W$ feature transformations for different architectures. (c-d): Evolution of the $RQ = \text{tr}(X^T \Delta X) / \|X\|$ for W and $2W$ as before. Experiments run on the Cora dataset for 50 random seeds.