

---

# ResolvNet: A Graph Convolutional Network with multi-scale Consistency

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 It is by now a well known fact in the graph learning community that the presence of  
2 bottlenecks severely limits the ability of graph neural networks to propagate infor-  
3 mation over long distances. What so far has not been appreciated is that, counter-  
4 intuitively, also the presence of strongly connected sub-graphs may severely restrict  
5 information flow in common architectures. Motivated by this observation, we  
6 introduce the concept of multi-scale consistency. At the node level this concept  
7 refers to the retention of a connected propagation graph even if connectivity varies  
8 over a given graph. At the graph-level, multi-scale consistency refers to the fact  
9 that distinct graphs describing the same object at different resolutions should be  
10 assigned similar feature vectors. As we show, both properties are not satisfied by  
11 popular graph neural network architectures. To remedy these shortcomings, we  
12 introduce ResolvNet, a flexible graph neural network based on the mathematical  
13 concept of resolvents. We rigorously establish its multi-scale consistency theoret-  
14 ically and verify it in extensive experiments on real world data: Here networks  
15 based on this ResolvNet architecture prove expressive; out-performing baselines  
16 significantly on many tasks; in- and outside the multi-scale setting.

## 17 1 Introduction

18 Learning on graphs has developed into a rich and complex field, providing spectacular results on  
19 problems as varied as protein design [28], traffic forecasting [23], particle physics [38], recommender  
20 systems [10] and traditional tasks such as node- and graph classification [43, 44].

21 Despite their successes, graph neural networks (GNNs) are still plagued by fundamental issues:  
22 Perhaps best known is the phenomenon of oversmoothing, capturing the fact that node-features  
23 generated by common GNN architectures become less informative as network depth increases  
24 [22, 27]. From the perspective of information flow however deeper networks would be preferable, as  
25 a  $K$  layer message passing network [13], may only facilitate information exchange between nodes  
26 that are at most  $K$ -edges apart – a phenomenon commonly referred to as under-reaching [1, 41].  
27 However, even if information *is* reachable within  $K$  edges, the structure of the graph might not be  
28 conducive to communicating it between distant nodes: If bottlenecks are present in the graph at  
29 hand, information from an exponentially growing receptive field needs to be squashed into fixed-size  
30 vectors to pass through the bottleneck. This oversquashing-phenomenon [1, 41] prevents common  
31 architectures from propagating messages between distant nodes without information loss in the  
32 presence of bottlenecks.

33 What has so far not been appreciated within the graph learning community is that – somewhat counter-  
34 intuitively – also the presence of strongly connected subgraphs severely restricts the information  
35 flow within popular graph neural network architectures; as we establish in this work. Motivated by

36 this observation, we consider the setting of multi-scale graphs and introduce, define and study the  
 37 corresponding problem of multi-scale consistency for graph neural networks:

38 Multi-scale graphs are graphs whose edges are distributed on (at least) two scales: One large scale  
 39 indicating strong connections within certain (connected) clusters, and one regular scale indicating a  
 40 weaker, regular connectivity outside these subgraphs. The lack of multi-scale consistency of common  
 41 architectures then arises as two sides of the same coin: At the node level, prominent GNNs are unable  
 42 to consistently integrate multiple connectivity scales into their propagation schemes: They essentially  
 43 only propagate information along edges corresponding to the largest scale. At the graph level, current  
 44 methods are not stable to variations in resolution scale: Two graphs describing the same underlying  
 45 object at different resolutions are assigned vastly different feature vectors.

46 **Contributions:** We introduce the concept of multi-scale consistency for GNNs and study its two  
 47 defining characteristics at the node- and graph levels. We establish that common GNN architectures  
 48 suffer from a lack of multi-scale consistency and – to remedy this shortcoming – propose the  
 49 **ResolvNet** architecture. This method is able to consistently integrate multiple connectivity scales  
 50 occurring within graphs. At the node level, this manifests as ResolvNet – in contrast to common  
 51 architectures – not being limited to propagating information via a severely disconnected effective  
 52 propagation scheme, when multiple scales are present within a given graph. At the graph-level, this  
 53 leads to ResolvNet provably and numerically verifiably assigning similar feature vectors to graphs  
 54 describing the same underlying object at varying resolution scales; a property which – to the best of  
 55 our knowledge – no other graph neural network has demonstrated.

## 56 2 Multi-Scale Graphs and Multi-Scale Consistency

### 57 2.1 Multi-Scale Graphs

58 We are interested in graphs with edges distributed on (at least) two scales: A large scale indicating  
 59 strong connections within certain clusters, and a regular scale indicating a weaker, regular connectivity  
 60 outside these subgraphs. Before giving a precise definition, we consider two instructive examples:

61 **Example I. Large Weights:** A two-scale geometry as outlined above, might e.g. arise within  
 62 weighted graphs discretizing underlying continuous spaces: Here, edge weights are typically deter-  
 63 mined by the inverse discretization length ( $w_{ij} \sim 1/d_{ij}$ ), which might vary over the graph [30, 31].  
 64 Strongly connected sub-graphs would then correspond to clusters of nodes that are spatially closely  
 65 co-located. Alternatively, such different scales can occur in social networks; e.g. if edge-weights  
 66 are set to number of exchanged messages. Nodes representing (groups of) close friends would then  
 67 typically be connected by stronger edges than nodes encoding mere acquaintances, which would  
 68 typically have exchanged fewer messages.

69 Given such a weighted graph, we partition its weighted adjacency matrix  $W = W_{\text{reg.}} + W_{\text{high}}$  into a  
 70 disjoint sum over a part  $W_{\text{reg.}}$  containing only regular edge-weights and part  $W_{\text{high}}$  containing only  
 71 large edge-weights. This decomposition induces two graph structures on the common node set  $\mathcal{G}$ : We  
 72 set  $G_{\text{reg.}} := (\mathcal{G}, W_{\text{reg.}})$  and  $G_{\text{high}} := (\mathcal{G}, W_{\text{high}})$  (c.f. also Fig. 1).

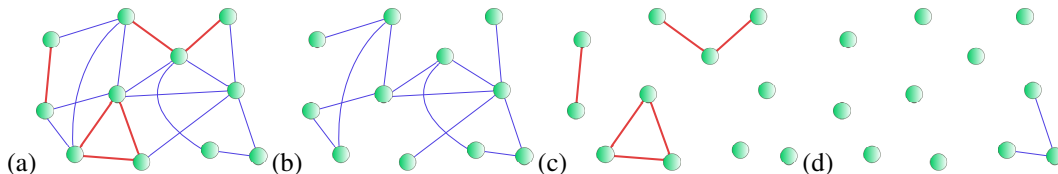


Figure 1: (a) Graph  $G$  with  $\mathcal{E}_{\text{reg.}}$  (blue) &  $\mathcal{E}_{\text{high}}$  (red); (b)  $G_{\text{reg.}}$ ; (c)  $G_{\text{high}}$ ; (d)  $G_{\text{excl.-reg.}}$ .

73 In preparation for our discussion in Section 2.2, we also define the graph  $G_{\text{excl.-reg.}}$  whose edges  
 74 consists of those elements  $(i, j) \in \mathcal{G} \times \mathcal{G}$  that do not have a neighbouring edge in  $G_{\text{high}}$ ; i.e. those  
 75 edges  $(i, j) \in \mathcal{E} \subseteq \mathcal{G} \times \mathcal{G}$  so that for any  $k \in \mathcal{G}$  we have  $(W_{\text{high}})_{ik}, (W_{\text{high}})_{kj} = 0$  (c.f. Fig. 1 (d)).

76 **Example 2. Many Connections:** Beyond weighted edges, disparate connectivities may also arise in  
 77 unweighted graphs with binary adjacency matrices: In a social network where edge weights encode a  
 78 binary friendship status for example, there might still exist closely knit communities within which  
 79 every user is friends with every other, while connections between such friend-groups may be sparser.

80 Here we may again split the adjacency matrix  $W = W_{\text{reg.}} + W_{\text{high}}$  into a disjoint sum over a part  
 81  $W_{\text{reg.}}$  encoding regular connectivity outside of tight friend groups and a summand  $W_{\text{high}}$  encoding  
 82 closely knit communities into dense matrix blocks. Fig. 2 depicts the corresponding graph structures.

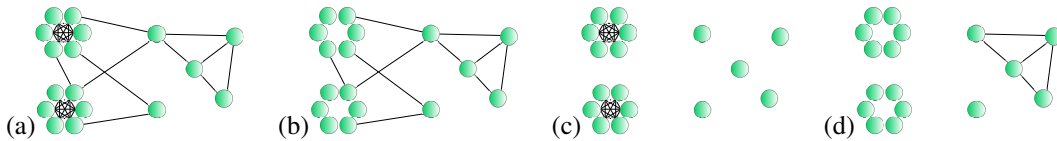


Figure 2: (a) Graph  $G$ ; (b)  $G_{\text{reg.}}$ ; (c)  $G_{\text{high}}$ ; (d)  $G_{\text{excl.-reg.}}$ .

83 **Exact Definition:** To unify both examples above into a common framework, we make use of tools  
 84 from spectral graph theory; namely the spectral properties of the **Graph Laplacian:** Given a graph  
 85  $G$  on  $N$  nodes, with weighted adjacency matrix  $W$ , diagonal degree matrix  $D$  and node weights  
 86  $\{\mu_i\}_{i=1}^N$  collected into the (diagonal) node-weight matrix  $M = \text{diag}(\{\mu_i\})$ , the (un-normalized)  
 87 graph Laplacian  $\Delta$  associated to the graph  $G$  is defined as  $\Delta = M^{-1}(D - W)$ .

88 It is a well known fact in spectral graph theory, that much information about the connectivity of the  
 89 graph  $G$  is encoded into the first (i.e. smallest) non-zero eigenvalue  $\lambda_1(\Delta)$  of this graph Laplacian  $\Delta$   
 90 [6, 7]. For an unweighted graph  $G$  on  $N$  nodes, this eigenvalue  $\lambda_1(\Delta)$  is for example maximised if  
 91 every node is connected to all other nodes (i.e.  $G$  is an  $N$ -clique); in which case we have  $\lambda_1(\Delta) = N$ .  
 92 For weighted graphs, it is clear that scaling all weights by a (large) constant  $c$  exactly also scales this  
 93 eigenvalue as  $\lambda_1(\Delta) \mapsto c \cdot \lambda_1(\Delta)$ . Thus the eigenvalue  $\lambda_1(\Delta)$  is indeed a good proxy for measuring  
 94 the strength of communities within a given graph  $G$ .

95 In order to formalize the concept of multi-scale graphs containing strongly connected subgraphs, we  
 96 thus make the following definition:

97 **Definition 2.1.** A Graph is called multi-scale if its weight-matrix  $W$  admits a *disjoint* decomposition

$$W = W_{\text{reg.}} + W_{\text{high}} \quad \text{with} \quad \lambda_1(\Delta_{\text{high}}) > \lambda_{\text{max}}(\Delta_{\text{reg.}}).$$

98 Note that this decomposition of  $W$  also implies  $\Delta = \Delta_{\text{reg.}} + \Delta_{\text{high}}$  for the respective Laplacians. Note  
 99 also that the graph-structure determined by  $G_{\text{high}}$  need not be completely connected for  $\lambda_1(\Delta_{\text{high}})$  to  
 100 be large (c.f. Fig.s 1 and 2 (c)): If there are multiple disconnected communities,  $\lambda_1(\Delta_{\text{high}})$  is given as  
 101 the minimal *non-zero* eigenvalue of  $\Delta_{\text{high}}$  restricted to these individual components of  $G_{\text{high}}$ . The  
 102 largest eigenvalue  $\lambda_{\text{max}}(\Delta_{\text{reg.}})$  of  $\Delta_{\text{reg.}}$  can be interpreted as measuring the "maximal connectivity"  
 103 within the graph structure  $G_{\text{reg.}}$ : By means of Gershgorin's circle theorem [2], we may bound it  
 104 as  $\lambda_{\text{max}}(\Delta_{\text{reg.}}) \leq 2 \cdot d_{\text{reg.,max}}$ , with  $d_{\text{reg.,max}}$  the maximal node-degree occurring in the graph  $G_{\text{reg.}}$ .  
 105 Hence  $\lambda_{\text{max}}(\Delta_{\text{reg.}})$  is small, if the connectivity within  $G_{\text{reg.}}$  is sparse.

## 106 2.2 Multi-Scale consistency

107 We are now especially interested in the setting where the scales occurring in a given graph  $G$  are well  
 108 separated (i.e.  $\lambda_1(\Delta_{\text{high}}) \gg \lambda_{\text{max}}(\Delta_{\text{reg.}})$ ). Below, we describe how graph neural networks should  
 109 ideally consistently incorporate such differing scales and detail how current architectures fail to do  
 110 so. As the influence of multiple scales within graphs manifests differently depending on whether  
 111 node-level- or graph-level tasks are considered, we will discuss these settings separately.

### 112 2.2.1 Node Level Consistency: Retention of connected propagation Graphs

113 The fundamental purpose of graph neural networks is that of generating node embeddings not only  
 114 dependent on local node-features, but also those of surrounding nodes. Even in the presence of  
 115 multiple scales in a graph  $G$ , it is thus very much desirable that information is propagated between  
 116 all nodes connected via the edges of  $G$  – and not, say, only along the dominant scale (i.e. via  $G_{\text{high}}$ ).

117 This is however not the case for popular graph neural network architectures: Consider for example  
 118 the graph convolutional network GCN [19]: Here, feature matrices  $X$  are updated via the update rule  
 119  $X \mapsto \hat{A} \cdot X$ , with the off-diagonal elements of  $\hat{A}$  given as  $\hat{A}_{ij} = W_{ij} / \sqrt{\hat{d}_i \cdot \hat{d}_j}$ . Hence the relative  
 120 importance  $\hat{A}_{ij}$  of a message between a node  $i$  of large (renormalised) degree  $\hat{d}_i \gg 1$  and a node  $j$   
 121 that is less strongly connected (e.g.  $\hat{d}_j = \mathcal{O}(1)$ ) is severely discounted.

122 In the presence of multiple scales as in Section 2.1, this thus leads to messages essentially only  
 123 being propagated over a disconnected effective propagation graph that is determined by the ef-  
 124 fective behaviour of  $\hat{A}$  in the presence of multiple scales. As we show in Appendix A using  
 the decompositions  $W = W_{\text{reg.}} + W_{\text{high}}$ , the matrix

125 
$$\hat{A} \approx \left( D_{\text{high}}^{-\frac{1}{2}} W_{\text{high}} D_{\text{high}}^{-\frac{1}{2}} + D_{\text{reg.}}^{-\frac{1}{2}} \tilde{W}_{\text{excl.-reg.}} D_{\text{reg.}}^{-\frac{1}{2}} \right)$$

Thus information is essentially only propagated within the connected components of  $G_{\text{high}}$  and via edges in  $G_{\text{excl.-reg.}}$  (detached from edges in  $G_{\text{high}}$ ).

126 Appendix A further details that this reduction to propagating information only along a disconnected  
 127 effective graph in the presence of multiple scales generically persists for popular methods (such as  
 128 e.g. attention based methods [42] or spectral methods [8]).

129 Propagating only over severely disconnected effective graphs as in Fig. 3 is clearly detrimental:

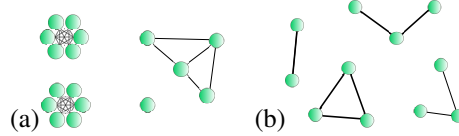


Figure 3: Effective propagation graphs for original graphs in Fig. 2 (a) and Fig. 1 (a)

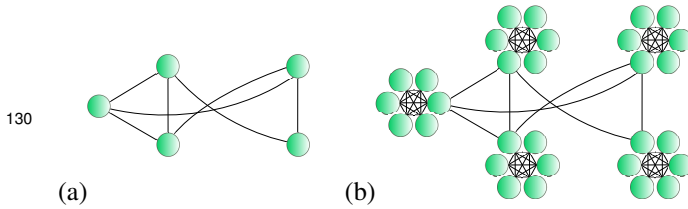


Figure 4: Individual nodes (a) replaced by 6-cliques (b)

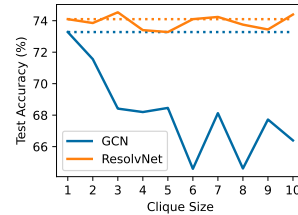


Figure 5: Classification Accuracy

131 As is evident from GCN’s performance in Fig.5, duplicating individual nodes of a popular graph  
 132 dataset into fully connected  $k$ -cliques as in Fig. 4 leads to a significant decrease in node-classification  
 133 accuracy, as propagation between cliques becomes increasingly difficult with growing clique-size  $k$ .  
 134 Details are provided in the Experimental-Section 5. In principle however, duplicating nodes does not  
 135 increase the complexity of the classification task at hand: Nodes and corresponding labels are only  
 136 duplicated in the train-, val.- and test-sets. What *is* changing however, is the geometry underlying the  
 137 problem; turning from a one-scale- into a two-scale setting with increasingly separated scales.

138 In Section 3 below, we introduce ResolvNet, which is able to consistently integrate multiple scales  
 139 within a given graph into its propagation scheme. As a result (c.f. Fig. 5) its classification accuracy is  
 140 not affected by an increasing clique-size  $k$  (i.e. an increasing imbalance in the underlying geometry).

## 141 2.2.2 Graph Level Consistency: Transferability between different Resolutions

142 At the graph level, we desire that graph-level feature vectors  $\Psi(G)$  generated by a network  $\Psi$  for  
 143 graphs  $G$  are stable to changes in resolution scales: More precisely, if two graphs  $G$  and  $\underline{G}$  describe  
 144 the same underlying object, space or phenomenon at different resolution scales, the generated feature  
 145 vectors should be close, as they encode *the same* object in the latent space. Ideally, we would  
 146 have a Lipschitz continuity relation that allows to bound the difference in generated feature vectors  
 147  $\|\Psi(G) - \Psi(\underline{G})\|$  in terms of a judiciously chosen distance  $d(G, \underline{G})$  between the graphs as

$$\|\Psi(G) - \Psi(\underline{G})\| \lesssim d(G, \underline{G}). \quad (1)$$

148 Note that a relation such as (1) also allows to make statements about *different* graphs  $G, \tilde{G}$  describing  
 149 an underlying object at *the same* resolution scale: If both such graphs are close to *the same* coarse  
 150 grained description  $\underline{G}$ , the triangle inequality yields  $\|\Psi(G) - \Psi(\tilde{G})\| \lesssim (d(G, \underline{G}) + d(\tilde{G}, \underline{G})) \ll 1$ .  
 151 To make precise what we mean by the coarse grained description  $\underline{G}$ , we revisit the example of  
 152 graphs discretising an underlying continuous space, with edge weights corresponding to inverse  
 153 discretization length ( $w_{ij} \sim 1/d_{ij}$ ). Coarse-graining – or equivalently lowering the resolution scale –  
 154 then corresponds to merging multiple spatially co-located nodes in the original graph  $G$  into single  
 155 aggregate nodes in  $\underline{G}$ . As distance scales inversely with edge-weight, this means that we are precisely  
 156 collapsing the strongly connected clusters within  $G_{\text{high}}$  into single nodes. Mathematically, we then  
 157 make this definition of the (lower resolution) coarse-grained graph  $\underline{G}$  exact as follows:

**Definition 2.2.** Denote by  $\underline{\mathcal{G}}$  the set of connected components in  $G_{\text{high}}$ . We give this set a graph structure  $\underline{G}$  as follows: Let  $R$  and  $P$  be elements of  $\underline{\mathcal{G}}$  (i.e. connected components in  $G_{\text{high}}$ ). We define the real number  $W_{RP}$  as  $W_{RP} = \sum_{r \in R} \sum_{p \in P} W_{rp}$ , with  $r$  and  $p$  nodes in the original graph  $G$ . We define the set of edges  $\mathcal{E}$  on  $\underline{G}$  as  $\mathcal{E} = \{(R, P) \in \underline{\mathcal{G}} \times \underline{\mathcal{G}} : W_{RP} > 0\}$  and assign  $W_{RP}$  as weight to such edges. Node weights of nodes in  $\underline{G}$  are defined similarly by aggregating weights of all nodes  $r$  contained in the connected component  $R$  of  $G_{\text{high}}$  as  $\mu_R = \sum_{r \in R} \mu_r$ .

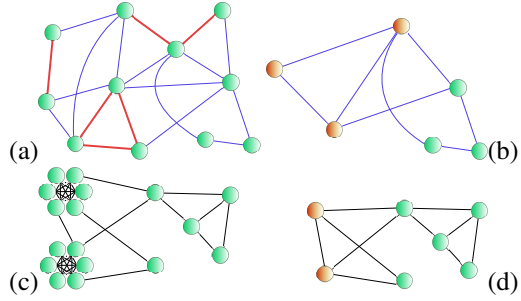


Figure 6: Original  $G$  (a,c) and coarsified  $\underline{G}$  (b,d)

This definition is of course also applicable to Example 2 of Section 2.1. Collapsing corresponding strongly connected component in a social network might then e.g. be interpreted as moving from interactions between individual users to considering interactions between (tightly-knit) communities.

While there have been theoretical investigations into this issue of **transferability** of graph neural networks between *distinct graphs* describing the *same system* [21, 34, 24, 20], the construction of an actual network with such properties – especially outside the asymptotic realm of very large graphs – has – to the best of our knowledge – so far not been successful. In Theorem 4.2 and Section 5 below, we show however that the ResolvNet architecture introduced in Section 3 below indeed provably and numerically verifiably satisfies (1), and is thus robust to variations in fine-print articulations of graphs describing the same object.

### 3 ResolvNet

We now design a network – termed ResolvNet – that can consistently incorporate multiple scales within a given graph into its propagation scheme. At the node level, we clearly want to avoid disconnected effective propagation schemes as discussed in Section 2.2.1 in settings with well-separated connectivity scales. At the graph level – following the discussion of Section 2.2.2 – we want to ensure that graphs  $G$  containing strongly connected clusters and graphs  $\underline{G}$  where these clusters are collapsed into single nodes are assigned similar feature vectors.

We can ensure both properties at the same time, if we manage to design a network whose propagation scheme when deployed on a multi-scale graph  $G$  is effectively described by propagating over a coarse grained version  $\underline{G}$  if the connectivity within the strongly connected clusters  $G_{\text{high}}$  of  $G$  is very large:

- At the node level, this avoids effectively propagating over disconnected limit graphs as in Section 2.2.1. Instead, information within strongly connected clusters is approximately homogenized and message passing is then performed on a (much better connected) coarse-grained version  $\underline{G}$  of the original graph  $G$  (c.f. Fig. 6).
- At the graph level, this means that the stronger the connectivity within the strongly connected clusters is, the more the employed propagation on  $G$  is like that on its coarse grained version  $\underline{G}$ . As we will see below, this can then be used to ensure the continuity property (1).

#### 3.1 The Resovent Operator

As we have seen in Section 2.2.1 (and as is further discussed in Appendix A), standard message passing schemes are unable to generate networks having our desired multi-scale consistency properties.

A convenient multi-scale description of graphs is instead provided by the graph Laplacian  $\Delta$  (c.f. Section 2.1), as this operator encodes information about coarse geometry of a graph  $G$  into small eigenvalues, while fine-print articulations of graphs correspond to large eigenvalues. [6, 7]. We are thus motivated to make use of this operator in our propagation scheme for ResolvNet.

In the setting of Example I of Section 2.1, letting the weights within  $G_{\text{high}}$  go to infinity (i.e. increasing the connectivity within the strongly connected clusters) however implies  $\|\Delta\| \rightarrow \infty$  for the norm of the Laplacian on  $G$ . Hence we *can not* implement propagation simply as  $X \mapsto \Delta \cdot X$ : This would not reproduce the corresponding propagation scheme on  $\underline{G}$  as we increase the connectivity within

197  $G_{\text{high}}$ : The Laplacian on  $G$  does not converge to the Laplacian on  $\underline{G}$  in the usual sense (it instead  
 198 diverges  $\|\Delta\| \rightarrow \infty$ ).

199 In order to capture convergence between operators with such (potentially) diverging norms, math-  
 200 ematicians have developed other – more refined – concepts: Instead of distances between original  
 201 operators, one considers distances between **resolvents** of such operators [40] :

202 **Definition 3.1.** The resolvent of an operator  $\Delta$  is defined as  $R_z(\Delta) := (\Delta - z \cdot Id)^{-1}$ , with  $Id$  the  
 203 identity mapping. Such resolvents are defined whenever  $z$  is not an eigenvalue of  $\Delta$ .

204 For Laplacians, taking  $z < 0$  hence ensures  $R_z(\Delta)$  is defined. Using this concept, we now rigorously  
 205 establish convergence (in the resolvent sense) of the Laplacian  $\Delta$  on  $G$  to the (coarse grained)  
 206 Laplacian  $\underline{\Delta}$  on  $\underline{G}$  as the connectivity within  $G_{\text{high}}$  is increased. To rigorously do so, we need to be  
 207 able to translate signals between the original graph  $G$  and its coarse-grained version  $\underline{G}$ :

208 **Definition 3.2.** Let  $x$  be a scalar graph signal. Denote by  $\mathbb{1}_R$  the vector that has 1 as entries on  
 209 nodes  $r$  belonging to the connected (in  $G_{\text{high}}$ ) component  $R$  and has entry zero for all nodes not in  
 210  $R$ . We define the down-projection operator  $J^\downarrow$  component-wise via evaluating at node  $R$  in  $\underline{G}$  as  
 211  $(J^\downarrow x)_R = \langle \mathbb{1}_R, x \rangle / \mu_R$ . This is then extended to feature *matrices*  $\{X\}$  via linearity. The interpolation  
 212 operator  $J^\uparrow$  is defined as  $J^\uparrow u = \sum_R u_R \cdot \mathbb{1}_R$ ; where  $u_R$  is a scalar value (the component entry of  $u$   
 213 at  $R \in \underline{G}$ ) and the sum is taken over all connected components of  $G_{\text{high}}$ .

214 With these preparations, we can rigorously establish that the *resolvent* of  $\Delta$  approaches that of  $\underline{\Delta}$ :

215 **Theorem 3.3.** We have  $R_z(\Delta) \rightarrow J^\uparrow R_z(\underline{\Delta}) J^\downarrow$  as connectivity within  $G_{\text{high}}$  increases. Explicitly:

$$\|R_z(\Delta) - J^\uparrow R_z(\underline{\Delta}) J^\downarrow\| = \mathcal{O}\left(\frac{\lambda_{\max}(\Delta_{\text{reg.}})}{\lambda_1(\Delta_{\text{high}})}\right)$$

216 The fairly involved proof of Theorem 3.3 is contained in Appendix B and builds on previous work:  
 217 We extend preliminary results in [20] by establishing *omni-directional* transferability (c.f. Theorem  
 218 4.1 below) and go beyond the toy-example of expanding a *single* node into a fixed and connected  
 219 sub-graph with pre-defined edge-weights.

220 The basic idea behind ResolvNet is then to (essentially) implement message passing as  $X \mapsto$   
 221  $R_z(\Delta) \cdot X$ . Building on Theorem 3.3, Section 4 below then makes precise how this rigorously enforces  
 222 multiscale-consistency as introduced in Section 2.2 in the corresponding ResolvNet architecture.

### 223 3.2 The ResolvNet Architecture

224 Building on Section 3.1, we now design filters for which feature propagation essentially proceeds  
 225 along the coarsified graph of Definition 2.2 as opposed to the disconnected effective graphs of Section  
 226 2.2.1, if multiple – well separated – edge-weight scales are present.

227 To this end, we note that Theorem 3.3 states for  $\lambda_1(\Delta_{\text{high}}) \gg \lambda_{\max}(\Delta_{\text{reg.}})$  (i.e. well separated scales),  
 228 that applying  $R_z(\Delta)$  to a signal  $x$  is essentially the same as first projecting  $x$  to  $\underline{G}$  via  $J^\downarrow$ , then  
 229 applying  $R_z(\underline{\Delta})$  there and finally lifting back to  $G$  with  $J^\uparrow$ . Theorem B.4 In Appendix B establishes  
 230 that this behaviour also persists for powers of resolvents; i.e. we also have  $R_z^k(\Delta) \approx J^\uparrow R_z^k(\underline{\Delta}) J^\downarrow$ .

231 **Resolvent filters:** This motivates us to choose our learnable filters as polynomials in resolvents

$$f_{z,\theta}(\Delta) := \sum_{k=a}^K \theta_k [(\Delta - zId)^{-1}]^k \quad (2)$$

232 with learnable parameters  $\{\theta_k\}_{k=a}^K$ . Thus our method can be interpreted as a spectral method [8],  
 233 with learned functions  $f_{z,\theta}(\lambda) = \sum_{k=a}^K \theta_k (\lambda - z)^{-k}$  applied to the operator  $\Delta$  determining our  
 234 convolutional filters. The parameter  $a$ , which determines the starting index of the sum in (2), may  
 235 either be set to  $a = 0$  (Type-0) or  $a = 1$  (Type-I). As we show in Theorem 4.1 below, this choice will  
 236 determine transferability properties of our models based on such filters.

237 Irrespectively, both Type-0 and Type-I filters are able to learn a wide array of functions, as the  
 238 following theorem (proved in Appendix C) shows:

239 **Theorem 3.4.** Fix  $\epsilon > 0$  and  $z < 0$ . For arbitrary functions  $g, h : [0, \infty] \rightarrow \mathbb{R}$  with  $\lim_{\lambda \rightarrow \infty} g(\lambda) =$   
 240 const. and  $\lim_{\lambda \rightarrow \infty} h(\lambda) = 0$ , there are filters  $f_{z,\theta}^0, f_{z,\theta}^I$  of Type-0 and Type-I respectively such that  
 241  $\|f_{z,\theta}^0 - g\|_\infty, \|f_{z,\theta}^I - h\|_\infty < \epsilon$ .

242 **The ResolvNet Layer:** Collecting resolvent filters into a convolutional architecture, the layer  
 243 wise update rule is then given as follows: Given a feature matrix  $X^\ell \in \mathbb{R}^{N \times F_\ell}$  in layer  $\ell$ , with  
 244 column vectors  $\{X_j^\ell\}_{j=1}^{F_\ell}$ , the feature vector  $X_i^{\ell+1}$  in layer  $\ell + 1$  is then calculated as  $X_i^{\ell+1} =$   
 245  $\text{ReLu}\left(\sum_{j=1}^{F_\ell} f_{z,\theta_{ij}^{\ell+1}}(\Delta) \cdot X_j^\ell + b_i^{\ell+1}\right)$  with a learnable bias vector  $b_i^{\ell+1}$ . Collecting biases into a  
 246 matrix  $B^{\ell+1} \in \mathbb{R}^{N \times F_{\ell+1}}$ , we can efficiently implement this using matrix-multiplications as

$$X^{\ell+1} = \text{ReLu}\left(\sum_{k=a}^K (T - \omega Id)^{-k} \cdot X^\ell \cdot W_k^{\ell+1} + B^{\ell+1}\right)$$

247 with weight matrices  $\{W_k^{\ell+1}\}$  in  $\mathbb{R}^{F_\ell \times F_{\ell+1}}$ . Biases are implemented as  $b_i = \beta_i \cdot \mathbb{1}_G$ , with  $\mathbb{1}_G$  the  
 248 vector of all ones on  $G$  and  $\beta_i \in \mathbb{R}$  learnable. This is done to ensure that the effective propagation on  
 249  $\underline{G}$  (if well separated scales are present in  $G$ ) is not disturbed by non-transferable bias terms on the  
 250 level of entire networks. This can be traced back to the fact that  $J^\downarrow \mathbb{1}_G = \mathbb{1}_{\underline{G}}$  and  $J^\uparrow \mathbb{1}_{\underline{G}} = \mathbb{1}_G$ . A  
 251 precise discussion of this matter is contained in Appendix D.

252 **Graph level feature aggregation:** As we will also consider the prediction of *graph-level* properties  
 253 in our experimental Section 5 below, we need to sensibly aggregate node-level features into graph-  
 254 level features on *node-weighted* graphs: As opposed to standard aggregation schemes (c.f. e.g. [45]),  
 255 we define an aggregation scheme  $\Psi$  that takes into account node weights and maps a feature matrix  
 256  $X \in \mathbb{R}^{N \times F}$  to a graph-level feature vector  $\Psi(X) \in \mathbb{R}^F$  via  $\Psi(X)_j = \sum_{i=1}^N |X_{ij}| \cdot \mu_i$ .

## 257 4 Multi-Scale consistency and Stability Guarantees

258 **Node Level:** We now establish rigorously that instead of propagating along disconnected effective  
 259 graphs (c.f. Fig. 3), ResolvNet instead propagates node features along the coarse-grained graphs of  
 260 Fig. 6 if multiple separated scales are present:

261 **Theorem 4.1.** Let  $\Phi$  and  $\underline{\Phi}$  be the maps associated to ResolvNets with the same learned weight  
 262 matrices and biases but deployed on graphs  $G$  and  $\underline{G}$  as defined in Section 3. We have

$$\|\Phi(J^\uparrow \underline{X}) - J^\uparrow \Phi(\underline{X})\|_2 \leq (C_1(\mathcal{W}) \cdot \|\underline{X}\|_2 + C_2(\mathcal{W}, \mathcal{B})) \cdot \|R_z(\Delta) - J^\uparrow R_z(\underline{\Delta}) J^\downarrow\|$$

263 if the network is based on Type-0 resolvent filters (c.f. Section 3). Additionally, we have

$$\|\Phi(X) - J^\downarrow \Phi(J^\downarrow X)\|_2 \leq (C_1(\mathcal{W}) \cdot \|X\|_2 + C_2(\mathcal{W}, \mathcal{B})) \cdot \|R_z(\Delta) - J^\downarrow R_z(\underline{\Delta}) J^\uparrow\|$$

264 if only Type-I filters are used in the network. Here  $C_1(\mathcal{W})$  and  $C_2(\mathcal{W}, \mathcal{B})$  are constants that depend  
 265 polynomially on singular values of learned weight matrices  $\mathcal{W}$  and biases  $\mathcal{B}$ .

266 The proof – as well as additional results – may be found in Appendix E. Note that Theorem 3.3  
 267 implies that both equations tends to zero for increasing scale separation  $\lambda_1(\Delta_{\text{high}}) \gg \lambda_{\max}(\Delta_{\text{reg}})$ .

268 The difference between utilizing Type-0 and Type-I resolvent filters, already alluded to in the  
 269 preceding Section 3, now can be understood as follows: Networks based on Type-0 filters effective-  
 270 ly propagate signals *lifted* from the coarse grained graph  $\underline{G}$  to the original graph  $G$  along  $\underline{G}$  if  
 271  $\lambda_1(\Delta_{\text{high}}) \gg \lambda_{\max}(\Delta_{\text{reg}})$ . In contrast – in the same setting – networks based on Type-I resolvent  
 272 filters effectively first *project* any input signal on  $G$  to  $\underline{G}$ , propagate there and then lift back to  $G$ .

273 **Graph Level:** Beyond a single graph, we also establish graph-level multi-scale consistency: As  
 274 discussed in Section 2.2.2, if two graphs describe the same underlying object (at different resolution  
 275 scales) corresponding feature vectors should be similar. This is captured by our next result:

276 **Theorem 4.2.** Denote by  $\Psi$  the aggregation method introduced in Section 3. With  $\mu(G) = \sum_{i=1}^N \mu_i$   
 277 the total weight of the graph  $G$ , we have in the setting of Theorem 4.1 with Type-I filters, that

$$\|\Psi(\Phi(X)) - \Psi(\underline{\Phi}(J^\downarrow X))\|_2 \leq \sqrt{\mu(G)} (C_1(\mathcal{W}) \|X\|_2 + C_2(\mathcal{W}, \mathcal{B})) \|R_z(\Delta) - J^\downarrow R_z(\underline{\Delta}) J^\uparrow\|.$$

278 This result thus indeed establishes the desired continuity relation (1), with the distance metric  $d(G, \underline{G})$   
 279 provided by the similarity  $\|R_z(\Delta) - J^\downarrow R_z(\underline{\Delta}) J^\uparrow\|$  of the resolvents of the two graphs.

280 **5 Experiments**

281 **Node Classification:** To establish that the proposed **ResolvNet** architecture **not only performs**  
 282 **well** in multi-scale settings, we conduct node classification experiments on multiple *un-weighted*  
 283 real world datasets, ranging in edge-homophily  $h$  from  $h = 0.11$  (very heterophilic), to  $h = 0.81$   
 284 (very homophilic). Baselines constitute an ample set of established and recent methods: Spectral  
 285 approaches, are represented by ChebNet [8], GCN [19], BernNet [15], ARMA [3] and MagNet [47].  
 286 Spatial methods are given by GAT [42], SAGE [14] and GIN [45]. We also consider PPNP [11] and  
 287 NSD [5]. Details on datasets, experimental setup and hyperparameters are provided in Appendix F.

Table 1: Average Accuracies [%] with uncertainties encoding the 95 % confidence Level. Top three models are coloured-coded as **First, Second, Third**.

$h$	MS. Acad. 0.81	Cora 0.81	Pubmed 0.80	Citeseer 0.74	Cornell 0.30	Actor 0.22	Squirrel 0.22	Texas 0.11
SAGE	<b>91.75</b> $\pm 0.09$	80.68 $\pm 0.30$	74.42 $\pm 0.42$	72.68 $\pm 0.32$	86.01 $\pm 0.72$	28.88 $\pm 0.32$	25.99 $\pm 0.28$	88.92 $\pm 0.73$
GIN	72.93 $\pm 1.94$	74.12 $\pm 1.21$	74.59 $\pm 0.45$	68.11 $\pm 0.69$	65.58 $\pm 1.23$	23.69 $\pm 0.28$	24.91 $\pm 0.58$	72.64 $\pm 1.19$
GAT	89.49 $\pm 0.15$	80.12 $\pm 0.33$	77.12 $\pm 0.41$	<b>73.20</b> $\pm 0.37$	74.39 $\pm 0.93$	24.55 $\pm 0.28$	<b>27.22</b> $\pm 0.31$	75.31 $\pm 1.09$
NSD	90.78 $\pm 0.13$	70.34 $\pm 0.47$	69.74 $\pm 0.50$	64.39 $\pm 0.50$	<b>87.78</b> $\pm 0.65$	27.62 $\pm 0.39$	24.96 $\pm 0.27$	<b>91.64</b> $\pm 0.62$
PPNP	91.22 $\pm 0.13$	<b>83.77</b> $\pm 0.27$	<b>78.42</b> $\pm 0.31$	<b>73.25</b> $\pm 0.37$	71.93 $\pm 0.84$	25.93 $\pm 0.35$	23.69 $\pm 0.43$	70.73 $\pm 1.27$
ChebNet	<b>91.62</b> $\pm 0.10$	78.70 $\pm 0.37$	73.63 $\pm 0.43$	72.10 $\pm 0.43$	85.99 $\pm 0.10$	<b>29.51</b> $\pm 0.31$	25.68 $\pm 0.28$	<b>91.01</b> $\pm 0.59$
GCN	90.81 $\pm 0.10$	81.49 $\pm 0.36$	76.60 $\pm 0.44$	71.34 $\pm 0.45$	73.35 $\pm 0.88$	24.60 $\pm 0.28$	<b>30.40</b> $\pm 0.40$	76.16 $\pm 1.12$
MagNet	87.23 $\pm 0.16$	76.50 $\pm 0.42$	68.23 $\pm 0.44$	70.92 $\pm 0.49$	<b>87.15</b> $\pm 0.66$	<b>30.50</b> $\pm 0.32$	23.54 $\pm 0.32$	<b>90.84</b> $\pm 0.54$
ARMA	88.97 $\pm 0.18$	81.24 $\pm 0.24$	76.28 $\pm 0.35$	70.64 $\pm 0.45$	83.68 $\pm 0.80$	24.40 $\pm 0.45$	22.72 $\pm 0.42$	87.41 $\pm 0.73$
BernNet	91.37 $\pm 0.14$	<b>83.26</b> $\pm 0.24$	<b>77.24</b> $\pm 0.37$	73.11 $\pm 0.34$	<b>87.14</b> $\pm 0.57$	28.90 $\pm 0.45$	22.86 $\pm 0.32$	89.81 $\pm 0.68$
ResolvNet	<b>92.73</b> $\pm 0.08$	<b>84.16</b> $\pm 0.26$	<b>79.29</b> $\pm 0.36$	<b>75.03</b> $\pm 0.29$	84.92 $\pm 1.43$	<b>29.06</b> $\pm 0.32$	<b>26.51</b> $\pm 0.23$	87.73 $\pm 0.89$

288 As is evident from Table 1, **ResolvNet out-performs all baselines in the homophilic setting**. This  
 289 can be traced back to the inductive bias ResolvNet is equipped with by design: It might be summarized  
 290 as "Nodes that are strongly connected should be assigned similar feature vectors" (c.f. Theorem 4.1) .  
 291 This inductive bias – necessary to achieve a consistent incorporation of multiple scales – is of course  
 292 counterproductive in the presence of heterophily; here nodes that are connected by edges generically  
 293 have *differing* labels and should thus be assigned different feature vectors. However the ResolvNet  
 294 architecture also performs well on most heterophilic graphs: It e.g. out-performs NSD – a recent  
 295 state of the art method specifically developed for heterophily – on two such graphs.

296 **Node Classification for increasingly separated scales:** To test ResolvNet’s ability to consistently  
 297 incorporate multiple scales in the unweighted setting against a representative baseline, we duplicated  
 298 individual nodes on the Citeseer dataset [36]  $k$ -times to form (fully connected)  $k$ -cliques; keeping  
 299 the train-val-test partition constant. GCN and ResolvNet were then trained on the same ( $k$ -fold  
 300 expanded) train-set and asked to classify nodes on the ( $k$ -fold expanded) test-partition. As discussed  
 301 in Section 1 (c.f. Fig.5) GCN’s performance decreased significantly, while ResolvNet’s accuracy  
 302 stayed essentially constant; showcasing its ability to consistently incorporate multiple scales.

303 **Regression on real-world multi-scale graphs:** In order to showcase the properties of our newly  
 304 developed method on real world data admitting a two-scale behaviour, we evaluate on the task of  
 305 molecular property prediction. While ResolvNet is not designed for this setting, this task still allows  
 306 to fairly compare its expressivity and stability properties against other non-specialized graph neural  
 307 networks [17]. Our dataset (QM7; [35]) contains descriptions of 7165 organic molecules; each  
 308 containing hydrogen and up to seven types of heavy atoms. A molecule is represented by its Coulomb  
 309 matrix, whose off-diagonal elements  $C_{ij} = Z_i Z_j / |\vec{x}_i - \vec{x}_j|$  correspond to the Coulomb repulsion  
 310 between atoms  $i$  and  $j$ . We treat  $C$  as a weighted adjacency matrix. Prediction target is the molecular  
 311 atomization energy, which – crucially – depends on long range interaction within molecules [46].  
 312 However, with edge-weights  $C_{ij}$  scaling as inverse distance, long range propagation of information is  
 313 scale-suppressed in the graph determined by  $C$ , when compared to the much larger weights between  
 314 closer atoms. We choose Type-I filters in ResolvNet, set node weights as atomic charge ( $\mu_i = Z_i$ )  
 315 and use one-hot encodings of atomic charges  $Z_i$  as node-wise input features.



As is evident from Table 2, our method produces significantly lower mean-absolute-errors (MAEs) than baselines of Table 1 deployable on weighted graphs. We attribute this to the fact that our model allows for long range information propagation within each molecule, as propagation along corresponding edges is suppressed for baselines but not for our model (c.f. Section 2.2.1). Appendix contains additional experiments on QM9 [32]; finding similar performance for (long-range dependent) energy targets.

Table 2: QM7-MAE

QM7	MAE [ <i>kcal/mol</i> ]
BernNet	113.57 $\pm$ 62.90
GCN	61.32 $\pm$ 1.62
ChebNet	59.57 $\pm$ 1.58
ARMA	59.39 $\pm$ 1.79
ResolvNet	<b>16.52<math>\pm</math>0.67</b>

**Stability to varying the resolution-scale:** To numerically verify the Stability-Theorem 4.2 – which guarantees similar graph-level feature vectors for graphs describing the same underlying object at different resolution scales – we conduct additional experiments: We modify (all) molecular graphs of QM7 by deflecting hydrogen atoms (H) out of their equilibrium positions towards the respective nearest heavy atom. This introduces a two-scale setting precisely as discussed in section 2: Edge weights between heavy atoms remain the same, while Coulomb repulsions between H-atoms and respective nearest heavy atom increasingly diverge. Given an original molecular graph  $G$  with node weights  $\mu_i = Z_i$ , the corresponding coarse-grained graph  $\underline{G}$  corresponds to a description where heavy atoms and surrounding H-atoms are aggregated into single super-nodes. Node-features of aggregated nodes are now no longer one-hot encoded charges, but normalized bag-of-word vectors whose individual entries encode how much of the total charge of a given super-node is contributed by individual atom-types. Appendix F provides additional details and examples.

In this setting, we now compare features generated for coarsified graphs  $\{\underline{G}\}$ , with feature generated for graphs  $\{G\}$  where hydrogen atoms have been deflected but have not yet completely arrived at the positions of nearest heavy atoms. Feature vectors are generated with the previously trained networks of Table 2. A corresponding plot is presented in Figure 7. Features generated by ResolvNet converge as the larger scale increases (i.e. the distance between hydrogen and heavy atoms decreases). This result numerically verifies the scale-invariance Theorem 4.2. As reference, we also plot the norm differences corresponding to baselines, which do not decay. We might thus conclude that these models – as opposed to ResolvNet – are scale- and resolution sensitive when generating graph level features. For BernNet we observe a divergence behaviour, which we attribute to numerical instabilities.

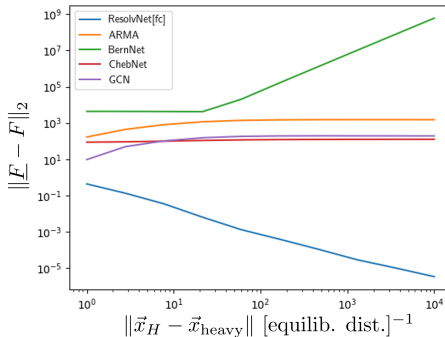


Figure 7: Feature-vector-difference for collapsed ( $\underline{F}$ ) and deformed ( $F$ ) graphs.

As a final experiment, we treat the coarse-grained molecular graphs  $\{\underline{G}\}$  as a model for data obtained from a resolution-limited observation process, that is unable to resolve positions of hydrogen and only provides information about how many H-atoms are bound to a given heavy atom. Given models trained on higher resolution data, atomization energies for such observed molecules are now to be predicted. Table 3 contains corresponding results. While the performance of baselines decreases significantly if the resolution scale is varied during inference, the prediction accuracy of ResolvNet remains high; even slightly increasing. While ResolvNet out-performed baselines by a factor of three on same-resolution-scale data (c.f. Table 2), its lead increases to a factor of 10 and higher in the multi-scale setting.

Table 3: QM7<sub>coarse</sub>-MAE

QM7	MAE [ <i>kcal/mol</i> ]
BernNet	580.67 $\pm$ 99.27
GCN	124.53 $\pm$ 34.58
ChebNet	645.14 $\pm$ 34.59
ARMA	248.96 $\pm$ 15.56
ResolvNet	<b>16.23<math>\pm</math>2.74</b>

## 6 Conclusion

This work introduced the concept of multi-scale consistency: At the node level this refers to the retention of a propagation scheme not solely determined by the largest given connectivity scale. At the graph-level it mandates that distinct graphs describing the same object at different resolutions should be assigned similar feature vectors. Common GNN architectures were shown to not be multi-scale consistent, while the newly introduced ResolvNet architecture was theoretically and experimentally established to have this property. Deployed on real world data, ResolvNet proved expressive and stable; out-performing baselines significantly on many tasks in- and outside the multi-scale setting.

342 **References**

- 343 [1] Uri Alon and Eran Yahav. On the bottleneck of graph neural networks and its practical  
344 implications. In *International Conference on Learning Representations*, 2021.
- 345 [2] Imre Bárány and József Solymosi. *Gershgorin Disks for Multiple Eigenvalues of Non-negative*  
346 *Matrices*, pages 123–133. Springer International Publishing, Cham, 2017.
- 347 [3] Filippo Maria Bianchi, Daniele Grattarola, Lorenzo Francesco Livi, and Cesare Alippi. Graph  
348 neural networks with convolutional arma filters. *IEEE Transactions on Pattern Analysis and*  
349 *Machine Intelligence*, 44:3496–3507, 2019.
- 350 [4] L. C. Blum and J.-L. Reymond. 970 million druglike small molecules for virtual screening in  
351 the chemical universe database GDB-13. *J. Am. Chem. Soc.*, 131:8732, 2009.
- 352 [5] Cristian Bodnar, Francesco Di Giovanni, Benjamin Paul Chamberlain, Pietro Liò, and  
353 Michael M. Bronstein. Neural sheaf diffusion: A topological perspective on heterophily  
354 and oversmoothing in gnns. *CoRR*, abs/2202.04579, 2022.
- 355 [6] Andries E. Brouwer and Willem H. Haemers. *Spectra of Graphs*. New York, NY, 2012.
- 356 [7] F. R. K. Chung. *Spectral Graph Theory*. American Mathematical Society, 1997.
- 357 [8] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks  
358 on graphs with fast localized spectral filtering. *Advances in neural information processing*  
359 *systems*, 29, 2016.
- 360 [9] Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric.  
361 In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- 362 [10] Chen Gao, Yu Zheng, Nian Li, Yinfeng Li, Yingrong Qin, Jinghua Piao, Yuhan Quan, Jianxin  
363 Chang, Depeng Jin, Xiangnan He, and Yong Li. A survey of graph neural networks for  
364 recommender systems: Challenges, methods, and directions, 2023.
- 365 [11] Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate:  
366 Graph neural networks meet personalized pagerank. In *7th International Conference on*  
367 *Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net,  
368 2019.
- 369 [12] Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate:  
370 Graph neural networks meet personalized pagerank. In *7th International Conference on*  
371 *Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net,  
372 2019.
- 373 [13] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural  
374 message passing for quantum chemistry. In *International conference on machine learning*,  
375 pages 1263–1272. PMLR, 2017.
- 376 [14] William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on  
377 large graphs. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob  
378 Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information*  
379 *Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017,*  
380 *December 4-9, 2017, Long Beach, CA, USA*, pages 1024–1034, 2017.
- 381 [15] Mingguo He, Zhewei Wei, Zengfeng Huang, and Hongteng Xu. Bernnet: Learning arbitrary  
382 graph spectral filters via bernstein approximation. In Marc’ Aurelio Ranzato, Alina Beygelzimer,  
383 Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural*  
384 *Information Processing Systems 34: Annual Conference on Neural Information Processing*  
385 *Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 14239–14251, 2021.
- 386 [16] Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.

- 387 [17] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele  
388 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.  
389 In Hugo Larochelle, Marc’ Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-  
390 Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference  
391 on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual,*  
392 2020.
- 393 [18] Tosio Kato. *Perturbation theory for linear operators; 2nd ed.* Grundlehren der mathematischen  
394 Wissenschaften : a series of comprehensive studies in mathematics. Springer, Berlin, 1976.
- 395 [19] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional  
396 networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon,  
397 France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- 398 [20] Christian Koke. Limitless stability for graph convolutional networks. In *11th International  
399 Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenRe-  
400 view.net, 2023.
- 401 [21] Ron Levie, Michael M. Bronstein, and Gitta Kutyniok. Transferability of spectral graph  
402 convolutional neural networks. *CoRR*, abs/1907.12972, 2019.
- 403 [22] Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional net-  
404 works for semi-supervised learning. In Sheila A. McIlraith and Kilian Q. Weinberger, editors,  
405 *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the  
406 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium  
407 on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA,  
408 February 2-7, 2018*, pages 3538–3545. AAAI Press, 2018.
- 409 [23] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural  
410 network: Data-driven traffic forecasting. In *6th International Conference on Learning Rep-  
411 resentations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track  
412 Proceedings*. OpenReview.net, 2018.
- 413 [24] Sohir Maskey, Ron Levie, and Gitta Kutyniok. Transferability of graph neural networks: an  
414 extended graphon approach. *CoRR*, abs/2109.10096, 2021.
- 415 [25] Andrew McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating the  
416 construction of internet portals with machine learning. *Inf. Retr.*, 3(2):127–163, 2000.
- 417 [26] Galileo Namata, Ben London, Lise Getoor, and Bert Huang. Query-driven active surveying for  
418 collective classification. 2012.
- 419 [27] Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for  
420 node classification. In *8th International Conference on Learning Representations, ICLR 2020,  
421 Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- 422 [28] et. al. Pablo Gainza. Deciphering interaction fingerprints from protein molecular surfaces using  
423 geometric deep learning. *Nature*, 2023.
- 424 [29] Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn:  
425 Geometric graph convolutional networks, 2020.
- 426 [30] Olaf. Post. *Spectral Analysis on Graph-like Spaces / by Olaf Post*. Lecture Notes in Mathematics,  
427 2039. Springer Berlin Heidelberg, Berlin, Heidelberg, 1st ed. 2012. edition, 2012.
- 428 [31] Olaf Post and Jan Simmer. Graph-like spaces approximated by discrete graphs and applications.  
429 *Mathematische Nachrichten*, 294(11):2237–2278, 2021.
- 430 [32] Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole von Lilienfeld.  
431 Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data*, 1, 2014.
- 432 [33] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding. *J.  
433 Complex Networks*, 9(2), 2021.

- 434 [34] Luana Ruiz, Luiz F. O. Chamon, and Alejandro Ribeiro. Graphon neural networks and the  
435 transferability of graph neural networks. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia  
436 Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information  
437 Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020,  
438 NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 439 [35] M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld. Fast and accurate modeling  
440 of molecular atomization energies with machine learning. *Physical Review Letters*, 108:058301,  
441 2012.
- 442 [36] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-  
443 Rad. Collective classification in network data. *AI Magazine*, 29(3):93, Sep. 2008.
- 444 [37] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann.  
445 Pitfalls of graph neural network evaluation. *CoRR*, abs/1811.05868, 2018.
- 446 [38] Jonathan Shlomi, Peter Battaglia, and Jean-Roch Vlimant. Graph neural networks in particle  
447 physics. *Machine Learning: Science and Technology*, 2(2):021001, jan 2021.
- 448 [39] Jie Tang, Jimeng Sun, Chi Wang, and Zi Yang. Social influence analysis in large-scale networks.  
449 In John F. Elder IV, Françoise Fogelman-Soulié, Peter A. Flach, and Mohammed Javeed  
450 Zaki, editors, *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge  
451 Discovery and Data Mining, Paris, France, June 28 - July 1, 2009*, pages 807–816. ACM, 2009.
- 452 [40] Gerald Teschl. *Mathematical Methods in Quantum Mechanics*. American Mathematical Society,  
453 2014.
- 454 [41] Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, and  
455 Michael M. Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature,  
456 2021.
- 457 [42] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua  
458 Bengio. Graph attention networks. In *6th International Conference on Learning Representations,  
459 ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*.  
460 OpenReview.net, 2018.
- 461 [43] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A  
462 comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and  
463 Learning Systems*, 32(1):4–24, jan 2021.
- 464 [44] Shunxin Xiao, Shiping Wang, Yuanfei Dai, and Wenzhong Guo. Graph neural networks in node  
465 classification: survey and evaluation. *Mach. Vis. Appl.*, 33(1):4, 2022.
- 466 [45] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
467 networks? In *7th International Conference on Learning Representations, ICLR 2019, New  
468 Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- 469 [46] Linfeng Zhang, Han Wang, Maria Carolina Muniz, Athanassios Z. Panagiotopoulos, Roberto  
470 Car, and Weinan E. A deep potential model with long-range electrostatic interactions. *The  
471 Journal of Chemical Physics*, 156(12), mar 2022.
- 472 [47] Xitong Zhang, Yixuan He, Nathan Brugnone, Michael Perlmutter, and Matthew J. Hirn. Magnet:  
473 A neural network for directed graphs. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N.  
474 Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information  
475 Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021,  
476 NeurIPS 2021, December 6-14, 2021, virtual*, pages 27003–27015, 2021.

477 **A Effective Propagation Schemes**

478 For definiteness, we here discuss limit-propagation schemes in the setting where **edge-weights** are  
 479 large. The discussion for high-connectivity in the Sense of Example II of Section 2.1 proceeds in  
 480 complete analogy.

481  
 482

483 In this section, we then take up again the setting of Section 2. We reformulate this setting here in  
 484 a slightly modified language, that is more adapted to discussing effective propagation schemes of  
 485 standard architectures:  
 486

487 We partition edges on a weighted graph  $G$ , into two disjoint sets  $\mathcal{E} = \mathcal{E}_{\text{reg.}} \dot{\cup} \mathcal{E}_{\text{high}}$ , where the set of  
 488 edges with large weights is given by:

$$\mathcal{E}_{\text{high}} := \{(i, j) \in \mathcal{E} : w_{ij} \geq S_{\text{high}}\}$$

489 and the set with small weights is given by:

$$\mathcal{E}_{\text{reg.}} := \{(i, j) \in \mathcal{E} : w_{ij} \leq S_{\text{reg.}}\}$$

490 for weight scales  $S_{\text{high}} > S_{\text{reg.}} > 0$ . Without loss of generality, assume  $S_{\text{reg.}}$  to be as low as possible  
 491 (i.e.  $S_{\text{reg.}} = \max_{(i,j) \in \mathcal{E}_{\text{reg.}}} w_{ij}$ ) and  $S_{\text{high}}$  to be as high as possible (i.e.  $S_{\text{high}} = \min_{(i,j) \in \mathcal{E}_{\text{high}}}$ ) and no  
 492 weights in between the scales.

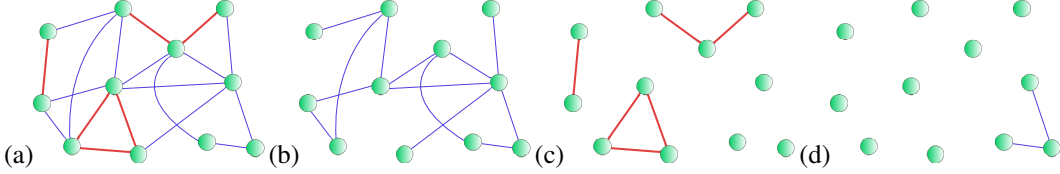


Figure 8: (a) Graph  $G$  with  $\mathcal{E}_{\text{reg.}}$  (blue) &  $\mathcal{E}_{\text{high}}$  (red); (b)  $G_{\text{reg.}}$ ; (c)  $G_{\text{high}}$ ; (d)  $G_{\text{reg., exclusive}}$

493 This decomposition induces two graph structures corresponding to the disjoint edge sets on the node  
 494 set  $\mathcal{G}$ : We set  $G_{\text{reg.}} := (\mathcal{G}, \mathcal{E}_{\text{reg.}})$  and  $G_{\text{high}} := (\mathcal{G}, \mathcal{E}_{\text{high}})$  c.f. Fig. 8).

495 We also introduce the set of edges  $\mathcal{E}_{\text{reg., exclusive}} := \{(i, j) \in \mathcal{E}_{\text{reg.}} \mid \forall k \in \mathcal{G} : (i, k) \notin \mathcal{E}_{\text{high}} \ \& \ (k, j) \notin \mathcal{E}_{\text{high}}\}$   
 496 connecting nodes that do not have an incident edge in  $\mathcal{E}_{\text{high}}$ . A corresponding example-graph  
 497  $G_{\text{reg., exclusive}}$  is depicted in Fig. 8 (d).

498

499 We are now interested in the behaviour of graph convolution schemes if the scales are well  
 500 separated:

$$S_{\text{high}} \gg S_{\text{reg.}}$$

501 **A.1 Spectral Convolutional Filters**

502 We first discuss resulting limit-propagation schemes for spectral convolutional networks. Such  
 503 networks implement convolutional filters as a mapping

$$x \longmapsto g_{\theta}(T)x$$

504 for a node feature  $x$ , a learnable function  $g_{\theta}$  and a graph shift operator  $T$ .

505 **A.1.1 Need for Normalization**

506 The graph shift operator  $T$  facilitating the graph convolutions needs to be normalized for established  
 507 spectral graph convolutional architectures:

508 For [3], this e.g. arises as a necessity for convergence of the proposed implementation scheme for the  
 509 rational filters introduced there (c.f. eq. (10) in [3]).

510 The work [8] needs its graph shift operator to be normalized, as it approximates generic filters  
 511 via a Chebyshev expansion. As argued in [8], such Chebyshev polynomials form an orthogonal  
 512 basis for the space  $L^2([-1, 1], dx/\sqrt{1-x^2})$ . Hence, the spectrum of the operator  $T$  to which the  
 513 (approximated and learned) function  $g_\theta$  is applied needs to be contained in the interval  $[-1, 1]$ .

514 In [19], it has been noted that for the architecture proposed there, choosing  $T$  to have eigenvalues in  
 515 the range  $[0, 2]$  (as opposed to the normalized ranges  $[0, 1]$  or  $[-1, 1]$ ) has the potential to lead to  
 516 vanishing- or exploding gradients as well as numerical instabilities. To alleviate this, [19] introduces  
 517 a "renormalization trick" (c.f. Section 2.2. of [19]) to produce a normalized graph shift operator on  
 518 which the network is then based.

519 We can understand the relationship between normalization of graph shift operator  $T$  and the stability  
 520 of corresponding convolutional filters explicitly: Assume that we have

$$\|T\| \gg 1.$$

521 This might e.g. happen when basing networks on the un-normalized graph Laplacian  $\Delta$  or the  
 522 weight-matrix  $W$  if edge weights are potentially large (such as in the setting  $S_{\text{high}} \gg S_{\text{reg}}$  that we are  
 523 considering).

524 By the spectral mapping theorem (see e.g. [40]), we have

$$\sigma(g_\theta(T)) = \{g_\theta(\lambda) : \lambda \in \sigma(T)\}, \quad (3)$$

525 with  $\sigma(T)$  denoting the spectrum (i.e. the set of eigenvalues) of  $T$ . For the largest (in absolute value)  
 526 eigenvalue  $\lambda_{\max}$  of  $T$ , we have

$$|\lambda_{\max}| = \|T\|. \quad (4)$$

527 Since learned functions are either implemented directly as a polynomial (as e.g. in [8, 15]) or  
 528 approximated as a Neumann type power iteration (as e.g. in [3, 12]) which can be thought of as a  
 529 polynomial, we have

$$\lim_{\lambda \rightarrow \pm\infty} |g_\theta(\lambda)| = \infty.$$

530 Thus in view of (3) and (4) we have for  $\|T\|$  sufficiently large, that

$$\|g_\theta(T)\| = |g_\theta(\pm\|T\|)|$$

531 with the sign  $\pm$  determined by  $\lambda_{\max} \geq 0$ . Since non-constant polynomials behave at least linearly  
 532 for large inputs, there is a constant  $C > 0$  such that

$$C \cdot \|T\| \leq \|g_\theta(T)\|$$

533 for all sufficiently large  $\|T\|$ . We thus have the estimate

$$\|x\| \cdot C \cdot \|T\| \leq \|g_\theta(T)x\|$$

534 for at least one input signal  $x$  (more precisely all  $x$  in the eigen-space corresponding to the largest (in  
 535 absolute value) eigenvalue  $\lambda_{\max}$ ). Thus if  $T$  is not normalized (i.e.  $\|T\|$  is not sufficiently bounded),  
 536 the norm of (hidden) features might increase drastically when moving from one (hidden) layer to the  
 537 next. This behaviour persists for all input signals  $x$  have components in eigenspaces corresponding to  
 538 large (in absolute value) eigenvalues of  $T$ .

### 539 A.1.2 Spectral Normalizations

As discussed in the previous Section A.1.1, instabilities arising from non-normalized graph shift operators can be traced back to the problem of such operators having large eigenvalues. It was thus – among other considerations – suggested in  
 540 [8] to base convolutional filters on the spectrally normalized graph shift operator

$$T = \frac{1}{\lambda_{\max}(\Delta)} \Delta,$$

541 with  $\Delta$  the un-normalized graph Laplacian. In the setting  $S_{\text{high}} \gg S_{\text{reg}}$  we are considering, this leads to  
 542 an effective feature propagation along  $G_{\text{high}}$  (c.f. also Fig. 9) only, as Theorem A.1 below establishes:  
 543

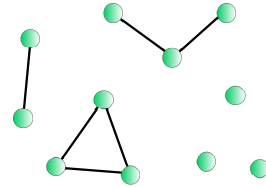


Figure 9: Limit graph corresponding to Fig 8 for spectral normalization

544 **Theorem A.1.** With

$$T = \frac{1}{\lambda_{\max}(\Delta)} \Delta,$$

545 and the scale decomposition as introduced in Section 2, we have that

$$\left\| T - \frac{1}{\lambda_{\max}(\Delta_{\text{high}})} \Delta_{\text{high}} \right\| = \mathcal{O}\left(\frac{S_{\text{reg.}}}{S_{\text{high}}}\right) \quad (5)$$

546 for  $S_{\text{high}} \gg S_{\text{reg.}}$ .

547 *Proof.* For convenience in notation, let us write

$$T_{\text{high}} = \frac{1}{\lambda_{\max}(\Delta_{\text{high}})} \Delta_{\text{high}}$$

548 and similarly

$$T_{\text{reg.}} = \frac{1}{\lambda_{\max}(\Delta_{\text{reg.}})} \Delta_{\text{reg.}}$$

549 In section 2, we already noted that

$$\Delta = \Delta_{\text{high}} + \Delta_{\text{reg.}},$$

550 which we may rewrite as

$$\Delta = \lambda_{\max}(\Delta_{\text{high}}) \cdot \left( T_{\text{high}} + \frac{\lambda_{\max}(\Delta_{\text{reg.}})}{\lambda_{\max}(\Delta_{\text{high}})} \cdot T_{\text{reg.}} \right). \quad (6)$$

551 Let us consider the equivalent expression

$$\frac{1}{\lambda_{\max}(\Delta_{\text{high}})} \cdot \Delta = T_{\text{high}} + \frac{\lambda_{\max}(\Delta_{\text{reg.}})}{\lambda_{\max}(\Delta_{\text{high}})} \cdot T_{\text{reg.}}. \quad (7)$$

552 We next note that

$$\lambda_{\max} \left( \frac{1}{\lambda_{\max}(\Delta_{\text{high}})} \cdot \Delta \right) = \frac{\lambda_{\max}(\Delta)}{\lambda_{\max}(\Delta_{\text{high}})}. \quad (8)$$

553 and

$$\lambda_{\max}(T_{\text{high}}) = 1$$

554 since the operation of taking eigenvalues of operators is multiplicative in the sense of

$$\lambda_{\max}(|a| \cdot T) = |a| \cdot \lambda_{\max}(T)$$

555 for non-negative  $|a| \geq 0$ .

556 Since the right-hand-side of (7) constitutes an analytic perturbation of  $T_{\text{high}}$ , we may apply analytic  
557 perturbation theory (c.f. e.g. [18] for an extensive discussion) to this problem. With this (together  
558 with  $\|T_{\text{high}}\| = 1$ ) we find

$$\lambda_{\max} \left( \frac{1}{\lambda_{\max}(\Delta_{\text{high}})} \cdot \Delta \right) = 1 + \mathcal{O} \left( \frac{\lambda_{\max}(\Delta_{\text{reg.}})}{\lambda_{\max}(\Delta_{\text{high}})} \right). \quad (9)$$

559 Using (8) and the fact that

$$\frac{\lambda_{\max}(\Delta_{\text{reg.}})}{\lambda_{\max}(\Delta_{\text{high}})} \propto \frac{S_{\text{reg.}}}{S_{\text{high}}}, \quad (10)$$

560 we thus have

$$\frac{\lambda_{\max}(\Delta)}{\lambda_{\max}(\Delta_{\text{high}})} = 1 + \mathcal{O} \left( \frac{S_{\text{reg.}}}{S_{\text{high}}} \right).$$

561 Since for small  $\epsilon$ , we also have

$$\frac{1}{1 + \epsilon} = 1 + \mathcal{O}(\epsilon),$$

562 the relation (10) also implies

$$\frac{\lambda_{\max}(\Delta_{\text{high}})}{\lambda_{\max}(\Delta)} = 1 + \mathcal{O} \left( \frac{S_{\text{reg.}}}{S_{\text{high}}} \right).$$

563 Multiplying (6) with  $1/\lambda_{\max}(\Delta)$  yields

$$T = \frac{\lambda_{\max}(\Delta_{\text{high}})}{\lambda_{\max}(\Delta)} \cdot \left( T_{\text{high}} + \frac{\lambda_{\max}(\Delta_{\text{reg.}})}{\lambda_{\max}(\Delta_{\text{high}})} \cdot T_{\text{reg.}} \right). \quad (11)$$

564 Since  $\|T_{\text{high}}\|, \|T_{\text{reg.}}\| = 1$  and

$$\frac{\lambda_{\max}(\Delta_{\text{reg.}})}{\lambda_{\max}(\Delta_{\text{high}})} \propto \frac{S_{\text{reg.}}}{S_{\text{high}}} < 1$$

565 for sufficiently large  $S_{\text{high}}$ , relation (11) implies

$$\left\| T - \frac{1}{\lambda_{\max}(\Delta_{\text{high}})} \Delta_{\text{high}} \right\| = \mathcal{O} \left( \frac{S_{\text{reg.}}}{S_{\text{high}}} \right)$$

566 as desired.

567 Note that we might in principle also make use of Lemma A.2 below, to provide quantitative bounds:

568 Lemma A.2 states that

$$|\lambda_k(A) - \lambda_k(B)| \leq \|A - B\|$$

569 for self-adjoint operators  $A$  and  $B$  and their respective  $k^{\text{th}}$  eigenvalues ordered by magnitude. On a  
 570 graph with  $N$  nodes, we clearly have  $\lambda_{\max} = \lambda_N$  for eigenvalues of (rescaled) graph Laplacians, since  
 571 all such eigenvalues are non-negative. This implies for the difference  $|1 - \lambda_{\max}(\Delta)/\lambda_{\max}(\Delta_{\text{high}})|$   
 572 arising in (9) that explicitly

$$\left| 1 - \frac{\lambda_{\max}(\Delta)}{\lambda_{\max}(\Delta_{\text{high}})} \right| \leq \frac{\lambda_{\max}(\Delta_{\text{reg.}})}{\lambda_{\max}(\Delta_{\text{high}})}.$$

573 This in turn can then be used to provide a quantitative bound in (5). Since we are only interested in  
 574 the qualitative behaviour for  $S_{\text{high}} \gg S_{\text{reg.}}$ , we shall however not pursue this further.

575 □

576 It remains to state and establish Lemma A.2 referenced at the end of the proof of Theorem A.1:

577

578 **Lemma A.2.** Let  $A$  and  $B$  be two hermitian  $n \times n$  dimensional matrices. Denote by  $\{\lambda_k(M)\}_{k=1}^n$   
 579 the eigenvalues of a hermitian matrix in increasing order.

580 With this we have:

$$|\lambda_k(A) - \lambda_k(B)| \leq \|A - B\|.$$

581 *Proof.* After the redefinition  $B \mapsto (-B)$ , what we need to prove is

$$|\lambda_i(A + B) - \lambda_i(A)| \leq \|B\|$$

582 for Hermitian  $A, B$ . Since we have

$$\lambda_i(A) - \lambda_i(A + B) = \lambda_i((A + B) + (-B)) - \lambda_i(A + B)$$

583 and  $\| -B \| = \|B\|$  it follows that it suffices to prove

$$\lambda_i(A + B) - \lambda_i(A) \leq \|B\|$$

584 for arbitrary hermitian  $A, B$ .

585 We note that the Courant-Fischer min – max theorem tells us that if  $A$  is an  $n \times n$  Hermitian matrix,  
 586 we have

$$\lambda_i(M) = \sup_{\dim(V)=i} \inf_{v \in V, \|v\|=1} v^* M v.$$



587 With this we find

$$\begin{aligned}
\lambda_i(A+B) - \lambda_i(A) &= \sup_{\dim(V)=i} \inf_{v \in V, \|v\|=1} v^*(A+B)v - \sup_{\dim(V)=i} \inf_{v \in V, \|v\|=1} v^*Av \\
&\leq \sup_{\dim(V)=i} \inf_{v \in V, \|v\|=1} v^*Av + \sup_{\dim(V)=i} \inf_{v \in V, \|v\|=1} v^*Bv \\
&\quad - \sup_{\dim(V)=i} \inf_{v \in V, \|v\|=1} v^*Av \\
&= \sup_{\dim(V)=i} \inf_{v \in V, \|v\|=1} v^*Bv \\
&= \sup_{\dim(V)=i} \inf_{v \in V, \|v\|=1} v^*Bv \\
&\leq \max_{1 \leq k \leq n} \{|\lambda_k(B)|\} \\
&= \|B\|.
\end{aligned}$$

588

□

### 589 A.1.3 Symmetric Normalizations

Most common spectral graph convolutional networks (such as e.g. [15, 3, 8]) base the learnable filters that they propose on the symmetrically normalized graph Laplacian

$$590 \quad \mathcal{L} = Id - D^{-\frac{1}{2}}WD^{-\frac{1}{2}}.$$

In the setting  $S_{\text{high}} \gg S_{\text{reg}}$  we are considering, this leads to an effective feature propagation along edges in  $\mathcal{E}_{\text{high}}$  and  $\mathcal{E}_{\text{low, exclusive}}$  (c.f. also Fig. 10) only, as Theorem A.3 below establishes:

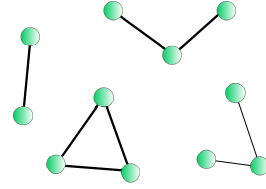


Figure 10: Limit graph corresponding to Fig 8 for symmetric normalization

591 **Theorem A.3.** With

$$T = Id - D^{-\frac{1}{2}}WD^{-\frac{1}{2}},$$

592 and the scale decomposition as introduced in Section 2, we have that

$$\left\| T - \left( Id - D_{\text{high}}^{-\frac{1}{2}}W_{\text{high}}D_{\text{high}}^{-\frac{1}{2}} - D_{\text{reg.}}^{-\frac{1}{2}}W_{\text{low, exclusive}}D_{\text{reg.}}^{-\frac{1}{2}} \right) \right\| = \mathcal{O} \left( \sqrt{\frac{S_{\text{reg.}}}{S_{\text{high}}}} \right) \quad (12)$$

593 for  $S_{\text{high}} \gg S_{\text{reg.}}$ .

594 *Proof.* We first note that instead of (12), we may equivalently establish

$$\left\| D^{-\frac{1}{2}}WD^{-\frac{1}{2}} - \left( D_{\text{high}}^{-\frac{1}{2}}W_{\text{high}}D_{\text{high}}^{-\frac{1}{2}} + D_{\text{reg.}}^{-\frac{1}{2}}W_{\text{low, exclusive}}D_{\text{reg.}}^{-\frac{1}{2}} \right) \right\| = \mathcal{O} \left( \sqrt{\frac{S_{\text{reg.}}}{S_{\text{high}}}} \right).$$

595 In Section 2, we already noted that

$$W = W_{\text{high}} + W_{\text{reg.}}$$

596 With this, we may write

$$D^{-\frac{1}{2}}WD^{-\frac{1}{2}} = D^{-\frac{1}{2}}W_{\text{high}}D^{-\frac{1}{2}} + D^{-\frac{1}{2}}W_{\text{reg.}}D^{-\frac{1}{2}}. \quad (13)$$

597 Let us first examine the term  $D^{-\frac{1}{2}}W_{\text{high}}D^{-\frac{1}{2}}$ . We note for the corresponding matrix entries that

$$\left( D^{-\frac{1}{2}}W_{\text{high}}D^{-\frac{1}{2}} \right)_{ij} = \frac{1}{\sqrt{d_i}} \cdot (W_{\text{high}})_{ij} \cdot \frac{1}{\sqrt{d_j}}$$

598 Let us use the notation

$$d_i^{\text{high}} = \sum_{j=1}^N (W_{\text{high}})_{ij}, \quad d_i^{\text{reg.}} = \sum_{j=1}^N (W_{\text{reg.}})_{ij} \quad \text{and} \quad d_i^{\text{low, exclusive}} = \sum_{j=1}^N (W_{\text{low, exclusive}})_{ij}.$$

599 We then find

$$\frac{1}{\sqrt{d_i}} = \frac{1}{\sqrt{d_i^{\text{high}}}} \cdot \frac{1}{\sqrt{1 + \frac{d_i^{\text{reg.}}}{d_i^{\text{high}}}}}$$

600 Using the Taylor expansion

$$\frac{1}{\sqrt{1 + \epsilon}} = 1 - \frac{1}{2}\epsilon + \mathcal{O}(\epsilon^2),$$

601 we thus have

$$\left(D^{-\frac{1}{2}}W_{\text{high}}D^{-\frac{1}{2}}\right)_{ij} = \frac{1}{\sqrt{d_i^{\text{high}}}} \cdot (W_{\text{high}})_{ij} \cdot \frac{1}{\sqrt{d_j^{\text{high}}}} + \mathcal{O}\left(\frac{d_i^{\text{reg.}}}{d_i^{\text{high}}}\right).$$

602 Since we have

$$\frac{d_i^{\text{reg.}}}{d_i^{\text{high}}} \propto \frac{S_{\text{reg.}}}{S_{\text{high}}},$$

603 this yields

$$D^{-\frac{1}{2}}W_{\text{high}}D^{-\frac{1}{2}} = D_{\text{high}}^{-\frac{1}{2}}W_{\text{high}}D_{\text{high}}^{-\frac{1}{2}} + \mathcal{O}\left(\frac{S_{\text{reg.}}}{S_{\text{high}}}\right).$$

604 Thus let us turn towards the second summand on the right-hand-side of (13). We have

$$\left(D^{-\frac{1}{2}}W_{\text{reg.}}D^{-\frac{1}{2}}\right)_{ij} = \frac{1}{\sqrt{d_i}} \cdot (W_{\text{reg.}})_{ij} \cdot \frac{1}{\sqrt{d_j}}.$$

605 Suppose that either  $i$  or  $j$  is not in  $G_{\text{low, exclusive}}$ . Without loss of generality (since the matrix under  
606 consideration is symmetric), assume  $i \notin G_{\text{low, exclusive}}$ , but  $(W_{\text{reg.}})_{ij} \neq 0$ . We may again write

$$\frac{1}{\sqrt{d_j}} = \frac{1}{\sqrt{d_j^{\text{high}}}} \cdot \frac{1}{\sqrt{1 + \frac{d_j^{\text{reg.}}}{d_j^{\text{high}}}}}.$$

607 Since

$$\frac{1}{\sqrt{1 + \frac{d_j^{\text{reg.}}}{d_j^{\text{high}}}}} \leq 1,$$

608 we have

$$\left|\left(D^{-\frac{1}{2}}W_{\text{reg.}}D^{-\frac{1}{2}}\right)_{ij}\right| \leq \left|\frac{1}{\sqrt{d_i}} \cdot (W_{\text{reg.}})_{ij}\right| \cdot \frac{1}{\sqrt{d_j^{\text{high}}}} = \mathcal{O}\left(\sqrt{\frac{S_{\text{reg.}}}{S_{\text{high}}}}\right).$$

609 If instead we have  $i, j \in G_{\text{low, exclusive}}$ , then clearly

$$\left(D^{-\frac{1}{2}}W_{\text{reg.}}D^{-\frac{1}{2}}\right)_{ij} = \left(D_{\text{reg.}}^{-\frac{1}{2}}W_{\text{low, exclusive}}D_{\text{reg.}}^{-\frac{1}{2}}\right)_{ij}.$$

610 Thus in total we have established

$$D^{-\frac{1}{2}}WD^{-\frac{1}{2}} = \left(D_{\text{high}}^{-\frac{1}{2}}W_{\text{high}}D_{\text{high}}^{-\frac{1}{2}} + D_{\text{reg.}}^{-\frac{1}{2}}W_{\text{low, exclusive}}D_{\text{reg.}}^{-\frac{1}{2}}\right) + \mathcal{O}\left(\frac{S_{\text{reg.}}}{S_{\text{high}}}\right)$$

611 which was to be established.

612 □

613 Apart from networks that make use of the symmetrically normalized graph Laplacian  $\mathcal{L}$ , some  
614 methods, such as most notably [19], instead base their filters on the operator

$$T = \tilde{D}^{-\frac{1}{2}}\tilde{W}\tilde{D}^{-\frac{1}{2}},$$

615 with

$$\tilde{W} = (W + Id)$$

616 and

$$\tilde{D} = D + Id.$$

617 In analogy to Theorem A.3, we here establish the limit propagation scheme determined by such  
618 operators:

619 **Theorem A.4.** With

$$T = \tilde{D}^{-\frac{1}{2}} \tilde{W} \tilde{D}^{-\frac{1}{2}},$$

620 where  $\tilde{W} = (W + Id)$  and  $\tilde{D} = D + Id$  as well as the scale decomposition of Section 2, we have that

$$\left\| T - \left( D_{\text{high}}^{-\frac{1}{2}} W_{\text{high}} D_{\text{high}}^{-\frac{1}{2}} + D_{\text{reg.}}^{-\frac{1}{2}} \tilde{W}_{\text{low, exclusive}} D_{\text{reg.}}^{-\frac{1}{2}} \right) \right\| = \mathcal{O} \left( \sqrt{\frac{S_{\text{reg.}} + 1}{S_{\text{high}}}} \right)$$

621 for  $S_{\text{high}} \gg S_{\text{reg.}}$ . Here  $\tilde{W}_{\text{low, exclusive}}$  is given as

$$\tilde{W}_{\text{low, exclusive}} := W_{\text{low, exclusive}} + \text{diag} \left( \mathbb{1}_{G_{\text{low, exclusive}}} \right)$$

622 and  $\mathbb{1}_{G_{\text{low, exclusive}}}$  denotes the vector whose entries are one for nodes in  $G_{\text{low, exclusive}}$  and zero for all  
623 other nodes.

624 The difference to the result of Theorem A.3 is thus that applicability of the limit propagation scheme  
625 of Fig. 10 for the GCN [19] is not only contingent upon  $S_{\text{high}} \gg S_{\text{reg.}}$  but also  $S_{\text{high}} \gg 1$ .

626 *Proof.* To establish this – as in the proof of Theorem A.3 – we first decompose  $T$ :

$$\begin{aligned} \tilde{D}^{-\frac{1}{2}} \tilde{W} \tilde{D}^{-\frac{1}{2}} &= \tilde{D}^{-\frac{1}{2}} W_{\text{high}} \tilde{D}^{-\frac{1}{2}} + \tilde{D}^{-\frac{1}{2}} W_{\text{reg.}} \tilde{D}^{-\frac{1}{2}} + \tilde{D}^{-\frac{1}{2}} Id \tilde{D}^{-\frac{1}{2}} \\ &= \tilde{D}^{-\frac{1}{2}} W_{\text{high}} \tilde{D}^{-\frac{1}{2}} + \tilde{D}^{-\frac{1}{2}} W_{\text{reg.}} \tilde{D}^{-\frac{1}{2}} + \tilde{D}^{-1} \end{aligned} \quad (14)$$

627 For the first term, we note

$$\left( \tilde{D}^{-\frac{1}{2}} W_{\text{high}} \tilde{D}^{-\frac{1}{2}} \right)_{ij} = \frac{1}{\sqrt{d_i + 1}} \cdot (W_{\text{high}})_{ij} \cdot \frac{1}{\sqrt{d_j + 1}}.$$

628 We then find

$$\frac{1}{\sqrt{d_i + 1}} = \frac{1}{\sqrt{d_i^{\text{high}}}} \cdot \frac{1}{\sqrt{1 + \frac{d_i^{\text{reg.}} + 1}{d_i^{\text{high}}}}}.$$

629 Analogously to the proof of Theorem A.3, this yields

$$\left( \tilde{D}^{-\frac{1}{2}} W_{\text{high}} \tilde{D}^{-\frac{1}{2}} \right)_{ij} = \frac{1}{\sqrt{d_i^{\text{high}}}} \cdot (W_{\text{high}})_{ij} \cdot \frac{1}{\sqrt{d_j^{\text{high}}}} + \mathcal{O} \left( \frac{1 + d_i^{\text{reg.}}}{d_i^{\text{high}}} \right).$$

630 This implies

$$\tilde{D}^{-\frac{1}{2}} W_{\text{high}} \tilde{D}^{-\frac{1}{2}} = D_{\text{high}}^{-\frac{1}{2}} W_{\text{high}} D_{\text{high}}^{-\frac{1}{2}} + \mathcal{O} \left( \frac{S_{\text{reg.}} + 1}{S_{\text{high}}} \right).$$

631 Next we turn to the second summand in (14):

$$\left( \tilde{D}^{-\frac{1}{2}} W_{\text{reg.}} \tilde{D}^{-\frac{1}{2}} \right)_{ij} = \frac{1}{\sqrt{d_i + 1}} \cdot (W_{\text{reg.}})_{ij} \cdot \frac{1}{\sqrt{d_j + 1}}.$$

632 Suppose that either  $i$  or  $j$  is not in  $G_{\text{low, exclusive}}$ . Without loss of generality (since the matrix under  
633 consideration is symmetric), assume  $i \notin G_{\text{low, exclusive}}$ , but  $(W_{\text{reg.}})_{ij} \neq 0$ . We may again write

$$\frac{1}{\sqrt{d_j + 1}} = \frac{1}{\sqrt{d_j^{\text{high}}}} \cdot \frac{1}{\sqrt{1 + \frac{d_i^{\text{reg.}} + 1}{d_i^{\text{high}}}}}.$$

634 Since

$$\frac{1}{\sqrt{1 + \frac{d_i^{\text{reg.}} + 1}{d_i^{\text{high}}}}} \leq 1,$$

635 we have

$$\begin{aligned}
\left| \left( D^{-\frac{1}{2}} W_{\text{reg.}} D^{-\frac{1}{2}} \right)_{ij} \right| &\leq \left| \frac{1}{\sqrt{1+d_i}} \cdot (W_{\text{reg.}})_{ij} \right| \cdot \frac{1}{\sqrt{d_j^{\text{high}}}} \\
&\leq \left| \frac{1}{\sqrt{d_i^{\text{reg.}}}} \cdot (W_{\text{reg.}})_{ij} \right| \cdot \frac{1}{\sqrt{d_j^{\text{high}}}} \\
&= \mathcal{O} \left( \sqrt{\frac{S_{\text{reg.}}}{S_{\text{high}}}} \right).
\end{aligned}$$

636 If instead we have  $i, j \in G_{\text{low, exclusive}}$ , then clearly

$$\left( \tilde{D}^{-\frac{1}{2}} W_{\text{reg.}} \tilde{D}^{-\frac{1}{2}} \right)_{ij} = \left( \tilde{D}_{\text{reg.}}^{-\frac{1}{2}} W_{\text{low, exclusive}} \tilde{D}_{\text{reg.}}^{-\frac{1}{2}} \right)_{ij}.$$

637 Finally we note for the third term on the right-hand-side of (14) that

$$\frac{1}{d_i} \leq \frac{1}{d_i^{\text{high}}} = \mathcal{O} \left( \frac{1}{S_{\text{high}}} \right)$$

638 if  $i \notin G_{\text{low, exclusive}}$ .

639 In total we thus have found

$$\tilde{D}^{-\frac{1}{2}} \tilde{W} \tilde{D}^{-\frac{1}{2}} = \left( D_{\text{high}}^{-\frac{1}{2}} W_{\text{high}} D_{\text{high}}^{-\frac{1}{2}} + D_{\text{reg.}}^{-\frac{1}{2}} \tilde{W}_{\text{low, exclusive}} D_{\text{reg.}}^{-\frac{1}{2}} \right) + \mathcal{O} \left( \sqrt{\frac{S_{\text{reg.}} + 1}{S_{\text{high}}}} \right);$$

640 which was to be proved. □

## 641 A.2 Spatial Convolutional Filters

642 Apart from spectral methods, there of course also exist methods that purely operate in the spatial  
643 domain of the graph. Such methods most often fall into the paradigm of message passing neural  
644 networks (MPNNs) [13, 9]: With  $X_i^\ell \in \mathbb{R}^F$  denoting the features of node  $i$  in layer  $\ell$  and  $w_{ij}$  denoting  
645 edge features, a message passing neural network may be described by the update rule (c.f. [13])

$$X_i^{\ell+1} = \gamma \left( X_i^\ell, \prod_{j \in \mathcal{N}(i)} \phi(X_i^\ell, X_j^\ell, w_{ij}) \right). \tag{15}$$

646 Here  $\mathcal{N}(i)$  denotes the neighbourhood of node  $i$ ,  $\prod$  denotes a differentiable and permutation invariant  
647 function (typically "sum", "mean" or "max") while  $\gamma$  and  $\phi$  denote differentiable functions such as  
648 multi-layer-perceptrons (MLPs) which might not be the same in each layer. [9].

649 Before we discuss corresponding limit-propagation schemes, we first establish that MPNNs are not  
650 able to reproduce the limit propagation scheme of Section 3 and are thus not stable to scale transitions  
651 and topological perturbations as discussed in Theorem 4.2 and Section 2.2.2.

### 652 A.2.1 Scale-Sensitivity of Message Passing Neural Networks

653 As we established in Theorem 4.1 and Theorem 4.2 (c.f. also the corresponding proofs in Appendix D  
654 and Appendix E respectively), the stability to scale-variations (such as coarse-graining) of ResolvNets  
655 arises from the reliance on *resolvents* and the limit propagation scheme that they establish if separated  
656 weight-scales are present (c.f. Appendix B below).

657 Here we establish that message passing networks (as defined in (15) above) are unable to emulate this  
658 limit propagation scheme. Hence such architectures are also not stable to scale-changing topological  
659 perturbations such as coarse-graining procedures.

To this end, we consider a simple, fully connected graph  $G$  on three nodes labeled 1, 2 and 3 (c.f. Fig. 11). We assume all node-weights to be equal to one ( $\mu_i = 1$  for  $i = 1, 2, 3$ ) and edge weights

$$w_{13}, w_{23} \leq S_{\text{reg.}}$$

as well as

$$w_{12} = S_{\text{high.}}$$

We now assume  $S_{\text{high}} \gg S_{\text{reg.}}$ .

Given states  $\{X_1^\ell, X_2^\ell, X_3^\ell\}$  in layer  $\ell$ , the limit propagation scheme introduced in Section 3 would require the updated feature vector of node 3 to be given by

$$X_{3,\text{desired}}^{\ell+1} := \gamma \left( X_3^\ell, \phi \left( X_3^\ell, \frac{X_1^\ell + X_2^\ell}{2}, (w_{31} + w_{32}) \right) \right)$$

However, the actual updated feature at node 3 is given as (c.f. (15)):

$$X_{3,\text{actual}}^{\ell+1} := \gamma \left( X_3^\ell, \phi \left( X_3^\ell, X_1^\ell, w_{31} \right) \llbracket \phi \left( X_3^\ell, X_2^\ell, w_{32} \right) \right) \quad (16)$$

Since there is no dependence on  $S_{\text{high}}$  in equation (16) – which defines  $X_{3,\text{actual}}^{\ell+1}$  – the desired propagation scheme can not arise, unless it is paradoxically already present at all scales  $S_{\text{high}}$ . If it is present at all scales, there is however only propagation along edges in  $G$ , even if  $S_{\text{high}} \approx S_{\text{reg.}}$ , which would imply that the message passing network would not respect the graph structure of  $G$ . Hence  $X_{3,\text{actual}}^{\ell+1} \rightarrow X_{3,\text{desired}}^{\ell+1}$  does not converge as  $S_{\text{high}}$  increases.

### A.2.2 Limit Propagation Schemes

The number of possible choices of message functions  $\phi$ , aggregation functions  $\llbracket \cdot \rrbracket$  and update functions  $\gamma$  is clearly endless. Here we shall exemplarily discuss limit propagation schemes for two popular architectures: We first discuss the most general case where the message function  $\phi$  is given as a learnable perceptron. Subsequently we assume that node features are updated with an attention-type mechanism.

**Generic message functions:** We first consider the possibility that the message function  $\phi$  in (16) is implemented via an MLP using ReLU-activations: Assuming (for simplicity in notation) a one-hidden-layer MLP mapping features  $X_i^\ell \in \mathbb{R}^{F_\ell}$  to features  $X_i^{\ell+1} \in \mathbb{R}^{F_{\ell+1}}$  we have

$$\phi(X_i^\ell, X_j^\ell, w_{ij}) = \text{ReLU} \left( W_1^\ell \cdot X_i^\ell + W_2^\ell \cdot X_j^\ell + W_3^\ell \cdot w_{ij} + B^\ell \right)$$

with bias term  $B^{\ell+1} \in \mathbb{R}^{F_{\ell+1}}$  and weight matrices  $W_1^{\ell+1}, W_2^{\ell+1} \in \mathbb{R}^{F_{\ell+1} \times F_\ell}$  and  $W_3^\ell \in \mathbb{R}^{F_{\ell+1}}$ .

We will assume that the weight-vector  $W_3^{\ell+1}$  has no-nonzero entries. This is not a severe limitation experimentally and in fact generically justified: The complementary event of at-least one entry of  $W_3$  being assigned precisely zero during training has probability weight zero (assuming an absolutely continuous probability distribution according to which weights are learned).

Let us now assume that the edge  $(ij)$  belongs to  $\mathcal{E}_{\text{high}}$  and the corresponding weight  $w_{ij}$  is large ( $w_{ij} \gg 1$ ). The behaviour of entries  $\phi(X_i^\ell, X_j^\ell, w_{ij})_a$  of the message  $\phi(X_i^\ell, X_j^\ell, w_{ij}) \in \mathbb{R}^{F_{\ell+1}}$  is then determined by the sign of the corresponding entry  $(W_3^\ell)_a$  of the weight vector  $W_3^\ell \in \mathbb{R}^{F_{\ell+1}}$ :

If we have  $(W_3^\ell)_a < 0$ , then  $\phi(X_i^\ell, X_j^\ell, w_{ij})_a$  approaches zero for larger edge-weights  $w_{ij}$ :

$$\lim_{w_{ij} \rightarrow \infty} \phi(X_i^\ell, X_j^\ell, w_{ij})_a = 0 \quad (17)$$

If we have  $(W_3^\ell)_a > 0$ , then  $\phi(X_i^\ell, X_j^\ell, w_{ij})_a$  increasingly diverges for larger edge-weights  $w_{ij}$ :

$$\lim_{w_{ij} \rightarrow \infty} \phi(X_i^\ell, X_j^\ell, w_{ij})_a = \infty \quad (18)$$

For either choice of aggregation function  $\llbracket \cdot \rrbracket$  in (15) among "max", "sum" or "mean" the behaviour in (18) leads to unstable networks if the update function  $\gamma$  is also given as an MLP with ReLU

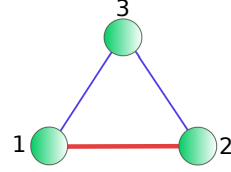


Figure 11: Three node Graph  $G$  with on large weight  $w_{12} \gg 1$ .

690 activations. Apart from instabilities, we also make the following observation: If  $S_{\text{high}} \gg S_{\text{reg}}$ , then by  
 691 (18) and continuity of  $\phi$  we can conclude that components  $\phi(X_i^\ell, X_j^\ell, w_{ij})_a$  of messages propagated  
 692 along  $\mathcal{E}_{\text{high}}$  for which  $(W_3^\ell)_a > 0$  dominate over messages propagated along edges in  $\mathcal{E}_{\text{reg}}$ . By (17),  
 693 the former clearly also dominate over components  $\phi(X_i^\ell, X_j^\ell, w_{ij})_a$  of messages propagated along  
 694  $\mathcal{E}_{\text{high}}$  for which  $(W_3^\ell)_a < 0$ . This behaviour is irrespective of whether "max", "sum" or "mean"  
 695 aggregations are employed. Hence the limit propagation scheme essentially only takes into account  
 696 message channels  $\phi(X_i^\ell, X_j^\ell, w_{ij})_a$  for which  $(ij) \in \mathcal{E}_{\text{high}}$  and  $(W_3^\ell)_a > 0$ .

697 Similar considerations apply, if non-linearities are chosen as leaky ReLU. If instead of ReLU  
 698 activations a sigmoid-nonlinearity  $\sigma$  like tanh is employed, messages propagated along  $\mathcal{E}_{\text{large}}$  become  
 699 increasingly uninformative, since they are progressively more independent of features  $X_i^\ell$  and weights  
 700  $w_{ij}$ . Indeed, for sigmoid activations, the limits (17) and (18) are given as follows:

701 If we have  $(W_3^\ell)_a < 0$ , then we have for larger edge-weights  $w_{ij}$  that

$$\lim_{w_{ij} \rightarrow \infty} \phi(X_i^\ell, X_j^\ell, w_{ij})_a = \lim_{y \rightarrow -\infty} \sigma(y).$$

702 If we have  $(W_3^\ell)_a > 0$ , then

$$\lim_{w_{ij} \rightarrow \infty} \phi(X_i^\ell, X_j^\ell, w_{ij})_a = \lim_{y \rightarrow \infty} \sigma(y).$$

703 In both cases, the messages  $\phi(X_i^\ell, X_j^\ell, w_{ij})$  propagated along  $\mathcal{E}_{\text{large}}$  become increasingly constant as  
 704 the scale  $S_{\text{high}}$  increases.

705 **Attention based messages:** Apart from general learnable message functions as above, we here  
 706 also discuss an approach where edge weights are re-learned in an attention based manner. For this we  
 707 modify the method [42] to include edge weights. The resulting propagation scheme – with a single  
 708 attention head for simplicity and a non-linearity  $\rho$  – is given as

$$X_i^{\ell+1} = \rho \left( \sum_{j \in \mathcal{N}(i)} \alpha_{ij} (W X_j^{\ell+1}) \right).$$

709 Here we have  $W \in \mathbb{R}^{F_{\ell+1} \times F_\ell}$  and

$$\alpha_{ij} = \frac{\exp(\text{LeakyRelu}(\vec{a}^\top [W X_i^\ell \parallel W X_j^\ell \parallel w_{ij}]))}{\sum_{k \in \mathcal{N}(i)} \exp(\text{LeakyRelu}(\vec{a}^\top [W X_i^\ell \parallel W X_k^\ell \parallel w_{ik}]))}, \quad (19)$$

710 with  $\parallel$  denoting concatenation. The weight vector  $\vec{a} \in \mathbb{R}^{2F_{\ell+1}+1}$  is assumed to have a non zero entry  
 711 in its last component. Otherwise, this attention mechanism would correspond to the one proposed  
 712 in [42], which does not take into account edge weights. Let us denote this entry of  $\vec{a}$  (determining  
 713 attention on the weight  $w_{ij}$ ) by  $a_w$ .

714 If  $a_w < 0$ , we have for  $(i, j) \in \mathcal{E}_{\text{high}}$  that

$$\exp(\text{LeakyRelu}(\vec{a}^\top [W X_i^\ell \parallel W X_j^\ell \parallel w_{ij}])) \rightarrow 0$$

715 as the weight  $w_{ij}$  increases. Thus propagation along edges in  $\mathcal{E}_{\text{high}}$  is essentially suppressed in this  
 716 case.

717 If  $a_w > 0$ , we have for  $(i, j) \in \mathcal{E}_{\text{high}}$  that

$$\exp(\text{LeakyRelu}(\vec{a}^\top [W X_i^\ell \parallel W X_j^\ell \parallel w_{ij}])) \rightarrow \infty$$

718 as the weight  $w_{ij}$  increases. Thus for edges  $(i, j) \in \mathcal{E}_{\text{reg}}$  (i.e. those that are *not* in  $\mathcal{E}_{\text{high}}$ ), we have

$$\alpha_{ij} \rightarrow 0,$$

719 since the denominator in (19) diverges. Hence in this case, propagation along  $\mathcal{E}_{\text{reg}}$  is essentially  
 720 suppressed and features are effectively only propagated along  $\mathcal{E}_{\text{high}}$ .

721 **B Proof of Theorem 3.3**

722 In this section, we prove Theorem 3.3. For convenience, we first restate the result – together with the  
723 definitions leading up to it – again:

724 **Definition B.1.** Denote by  $\underline{\mathcal{G}}$  the set of connected components in  $G_{\text{high}}$ . We give this set a graph  
725 structure as follows: Let  $R$  and  $P$  be elements of  $\underline{\mathcal{G}}$  (i.e. connected components in  $G_{\text{high}}$ ). We define  
726 the real number

$$W_{RP} = \sum_{r \in R} \sum_{p \in P} W_{rp},$$

727 with  $r$  and  $p$  nodes in the original graph  $G$ . We define the set of edges  $\underline{\mathcal{E}}$  on  $\underline{\mathcal{G}}$  as

$$\underline{\mathcal{E}} = \{(R, P) \in \underline{\mathcal{G}} \times \underline{\mathcal{G}} : W_{RP} > 0\}$$

728 and assign  $W_{RP}$  as weight to such edges. Node weights of limit nodes are defined similarly as  
729 aggregated weights of all nodes  $r$  (in  $G$ ) contained in the component  $R$  as

$$\underline{\mu}_R = \sum_{r \in R} \mu_r.$$

730 In order to translate signals between the original graph  $G$  and the limit description  $\underline{\mathcal{G}}$ , we need  
731 translation operators mapping signals from one graph to the other:

732 **Definition B.2.** Denote by  $\mathbb{1}_R$  the vector that has 1 as entries on nodes  $r$  belonging to the connected  
733 (in  $G_{\text{high}}$ ) component  $R$  and has entry zero for all nodes not in  $R$ . We define the down-projection  
734 operator  $J^\downarrow$  component-wise via evaluating at node  $R$  in  $\underline{\mathcal{G}}$  as

$$(J^\downarrow x)_R = \langle \mathbb{1}_R, x \rangle / \underline{\mu}_R.$$

735 The upsampling operator  $J^\uparrow$  is defined as

$$J^\uparrow u = \sum_R u_R \cdot \mathbb{1}_R; \tag{20}$$

736 where  $u_R$  is a scalar value (the component entry of  $u$  at  $R \in \underline{\mathcal{G}}$ ) and the sum is taken over all connected  
737 components in  $G_{\text{high}}$ .

738 The result we then have to prove is the following:

739 **Theorem B.3.** We have

$$\|R_z(\Delta) - J^\uparrow R_z(\underline{\Delta}) J^\downarrow\| = \mathcal{O}\left(\frac{\|\Delta_{\text{reg.}}\|}{\lambda_1(\Delta_{\text{high}})}\right)$$

740 holds; with  $\lambda_1(\Delta_{\text{high}})$  denoting the first non-zero eigenvalue of  $\Delta_{\text{high}}$ .

741 Note that this then indeed proves Theorem 3.3, since we have

$$\lambda_{\max}(\Delta_{\text{reg.}}) = \|\Delta_{\text{reg.}}\|.$$

742 *Proof.* We will split the proof of this result into multiple steps. For  $z < 0$  Let us denote by

$$\begin{aligned} R_z(\Delta) &= (\Delta - zId)^{-1}, \\ R_z(\Delta_{\text{high}}) &= (\Delta_{\text{high}} - zId)^{-1} \\ R_z(\Delta_{\text{reg.}}) &= (\Delta_{\text{reg.}} - zId)^{-1} \end{aligned}$$

743 the resolvents corresponding to  $\Delta$ ,  $\Delta_{\text{high}}$  and  $\Delta_{\text{reg.}}$  respectively.

744 Our first goal is establishing that we may write

$$R_z(\Delta) = [Id + R_z(\Delta_{\text{high}})\Delta_{\text{reg.}}]^{-1} \cdot R_z(\Delta_{\text{high}})$$

745 This will follow as a consequence of what is called the second resolvent formula [40]:

746 "Given self-adjoint operators  $A, B$ , we may write

$$R_z(A + B) - R_z(A) = -R_z(A)BR_z(A + B)."$$

747 In our case, this translates to

$$R_z(\Delta) - R_z(\Delta_{high}) = -R_z(\Delta_{high})\Delta_{reg}R_z(\Delta)$$

748 or equivalently

$$[Id + R_z(\Delta_{high})\Delta_{reg.}] R_z(\Delta) = R_z(\Delta_{high}).$$

749 Multiplying with  $[Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1}$  from the left then yields

$$R_z(\Delta) = [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} \cdot R_z(\Delta_{high})$$

750 as desired.

751 Hence we need to establish that  $[Id + R_z(\Delta_{high})\Delta_{reg.}]$  is invertible for  $z < 0$ .

752

753 To establish a contradiction, assume it is not invertible. Then there is a signal  $x$  such that

$$[Id + R_z(\Delta_{high})\Delta_{reg.}] x = 0.$$

754 Multiplying with  $(\Delta_{high} - zId)$  from the left yields

$$(\Delta_{high} + \Delta_{reg.} - zId)x = 0$$

755 which is precisely to say that

$$(\Delta - zId)x = 0$$

756 But since  $\Delta$  is a graph Laplacian, it only has non-negative eigenvalues. Hence we have reached our  
757 contradiction and established

$$R_z(\Delta) = [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} R_z(\Delta_{high}).$$

758

759 Our next step is to establish that

$$R_z(\Delta_{high}) \rightarrow \frac{P_0^{high}}{-z},$$

760 where  $P_0^{high}$  is the spectral projection onto the eigenspace corresponding to the lowest lying eigenvalue  
761  $\lambda_0(\Delta_{high}) = 0$  of  $\Delta_{high}$ . Indeed, by the spectral theorem for finite dimensional operators (c.f. e.g.  
762 [40]), we may write

$$R_z(\Delta_{high}) \equiv (\Delta_{high} - zId)^{-1} = \sum_{\lambda \in \sigma(\Delta_{high})} \frac{1}{\lambda - z} \cdot P_\lambda^{high}.$$

763 Here  $\sigma(\Delta_{high})$  denotes the spectrum (i.e. the collection of eigenvalues) of  $\Delta_{high}$  and the  
764  $\{P_\lambda^{high}\}_{\lambda \in \sigma(\Delta_{high})}$  are the corresponding (orthogonal) eigenprojections onto the eigenspaces of the  
765 respective eigenvalues. Thus we find

$$\left\| R_z(\Delta_{high}) - \frac{P_0^{high}}{-z} \right\| = \left\| \sum_{0 < \lambda \in \sigma(\Delta_{high})} \frac{1}{\lambda - z} \cdot P_\lambda^{high} \right\|;$$

766 where the sum on the right hand side now excludes the eigenvalue  $\lambda = 0$ .

767 Using orthonormality of the spectral projections, the fact that  $z < 0$  and monotonicity of  $1/(\cdot + |z|)$   
768 we find

$$\left\| R_z(\Delta_{high}) - \frac{P_0^{high}}{-z} \right\| = \frac{1}{\lambda_1(\Delta_{high}) + |z|}.$$

769 Here  $\lambda_1(\Delta_{high})$  is the first non-zero eigenvalue of  $(\Delta_{high})$ .

770 Non-zero eigenvalues scale linearly with the weight scale since we have

$$\lambda(S \cdot \Delta) = S \cdot \lambda(\Delta)$$

771 for any graph Laplacian (in fact any matrix)  $\Delta$  with eigenvalue  $\lambda$ . Thus we have

$$\left\| R_z(\Delta_{high}) - \frac{P_0^{high}}{-z} \right\| = \frac{1}{\lambda_1(\Delta_{high}) + |z|} \leq \frac{1}{\lambda_1(\Delta_{high})} \rightarrow 0$$



772 as  $\lambda_1(\Delta_{high}) \rightarrow \infty$ .

773

774 Our next task is to use this result in order to bound the difference

$$I := \left\| \left[ Id + \frac{P_0^{high}}{-z} \Delta_{reg.} \right]^{-1} \frac{P_0^{high}}{-z} - [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} R_z(\Delta_{high}) \right\|.$$

775 To this end we first note that the relation

$$[A + B - zId]^{-1} = [Id + R_z(A)B]^{-1} R_z(A)$$

776 provided to us by the second resolvent formula, implies

$$[Id + R_z(A)B]^{-1} = Id - B[A + B - zId]^{-1}.$$

777 Thus we have

$$\begin{aligned} \left\| [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} \right\| &\leq 1 + \|\Delta_{reg.}\| \cdot \|R_z(\Delta)\| \\ &\leq 1 + \frac{\|\Delta_{reg.}\|}{|z|}. \end{aligned}$$

778 With this, we have

$$\begin{aligned} &\left\| \left[ Id + \frac{P_0^{high}}{-z} \Delta_{reg.} \right]^{-1} \cdot \frac{P_0^{high}}{-z} - R_z(\Delta) \right\| \\ &= \left\| \left[ Id + \frac{P_0^{high}}{-z} \Delta_{reg.} \right]^{-1} \cdot \frac{P_0^{high}}{-z} - [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} \cdot R_z(\Delta_{high}) \right\| \\ &\leq \left\| \frac{P_0^{high}}{-z} \right\| \cdot \left\| \left[ Id + \frac{P_0^{high}}{-z} \Delta_{reg.} \right]^{-1} - [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} \right\| + \left\| \frac{P_0^{high}}{-z} - R_z(\Delta_{high}) \right\| \cdot \left\| [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} \right\| \\ &\leq \frac{1}{|z|} \left\| \left[ Id + \frac{P_0^{high}}{-z} \Delta_{reg.} \right]^{-1} - [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} \right\| + \left( 1 + \frac{\|\Delta_{reg.}\|}{|z|} \right) \cdot \frac{1}{\lambda_1(\Delta_{high})}. \end{aligned}$$

779 Hence it remains to bound the left hand summand. For this we use the following fact (c.f. [16],

780 Section 5.8. "Condition numbers: inverses and linear systems"):

781

782 Given square matrices  $A, B, C$  with  $C = B - A$  and  $\|A^{-1}C\| < 1$ , we have

$$\|A^{-1} - B^{-1}\| \leq \frac{\|A^{-1}\| \cdot \|A^{-1}C\|}{1 - \|A^{-1}C\|}.$$

783 In our case, this yields (together with  $\|P_0^{high}\| = 1$ ) that

$$\begin{aligned} &\left\| \left[ Id + P_0^{high}/(-z) \cdot \Delta_{reg.} \right]^{-1} - [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} \right\| \\ &\leq \frac{(1 + \|\Delta_{reg.}\|/|z|)^2 \cdot \|\Delta_{reg.}\| \cdot \left\| \frac{P_0^{high}}{-z} - R_z(\Delta_{high}) \right\|}{1 - (1 + \|\Delta_{reg.}\|/|z|) \cdot \|\Delta_{reg.}\| \cdot \left\| \frac{P_0^{high}}{-z} - R_z(\Delta_{high}) \right\|} \end{aligned}$$

784 For  $S_{high}$  sufficiently large, we have

$$\left\| -P_0^{high}/z - R_z(\Delta_{high}) \right\| \leq \frac{1}{2(1 + \|\Delta_{reg.}\|/|z|)}$$

785 so that we may estimate

$$\begin{aligned}
& \left\| \left[ Id + \Delta_{reg.} \frac{P_0^{high}}{-z} \right]^{-1} - [Id + \Delta_{reg.} R_z(\Delta_{high})]^{-1} \right\| \\
& \leq 2 \cdot (1 + \|\Delta_{reg.}\|) \cdot \left\| \frac{P_0^{high}}{-z} - R_z(\Delta_{high}) \right\| \\
& = 2 \frac{1 + \|\Delta_{reg.}\|/|z|}{\lambda_1(\Delta_{high})}
\end{aligned}$$

786 Thus we have now established

$$\left| \left[ Id + \frac{P_0^{high}}{-z} \Delta_{reg.} \right]^{-1} \cdot \frac{P_0^{high}}{-z} - R_z(\Delta) \right| = \mathcal{O} \left( \frac{\|\Delta_{reg.}\|}{\lambda_1(\Delta_{high})} \right).$$

787

788 Hence we are done with the proof, as soon as we can establish

$$\left[ -z Id + P_0^{high} \Delta_{reg.} \right]^{-1} P_0^{high} = J^\uparrow R_z(\underline{\Delta}) J^\downarrow,$$

789 with  $J^\uparrow, \underline{\Delta}, J^\downarrow$  as defined above. To this end, we first note that

$$J^\uparrow \cdot J^\downarrow = P_0^{high} \tag{21}$$

790 and

$$J^\downarrow \cdot J^\uparrow = Id_{\underline{G}}. \tag{22}$$

791 Indeed, the relation (21) follows from the fact that the eigenspace corresponding to the eigenvalue  
792 zero is spanned by the vectors  $\{\mathbb{1}_R\}_R$ , with  $\{R\}$  the connected components of  $G_{high}$ . Equation (22)  
793 follows from the fact that

$$\langle \mathbb{1}_R, \mathbb{1}_R \rangle = \underline{\mu}_R.$$

794 With this we have

$$\left[ Id + P_0^{high} \Delta_{reg.} \right]^{-1} P_0^{high} = [Id + J^\uparrow J^\downarrow \Delta_{reg.}]^{-1} J^\uparrow J^\downarrow.$$

795 To proceed, set

$$\underline{x} := F^\downarrow x$$

796 and

$$\mathcal{X} = \left[ P_0^{high} \Delta_{reg.} - z Id \right]^{-1} P_0^{high} x.$$

797 Then

$$\left[ P_0^{high} \Delta_{reg.} - z Id \right] \mathcal{X} = P_0^{high} x$$

798 and hence  $\mathcal{X} \in \text{Ran}(P_0^{high})$ . Thus we have

$$J^\uparrow J^\downarrow (\Delta_{reg.} - z Id) J^\uparrow J^\downarrow \mathcal{X} = J^\uparrow J^\downarrow x.$$

799 Multiplying with  $J^\downarrow$  from the left yields

$$J^\downarrow (\Delta_{reg.} - z Id) J^\uparrow J^\downarrow \mathcal{X} = J^\downarrow x.$$

800 Thus we have

$$(J^\downarrow \Delta_{reg.} J^\uparrow - z Id) J^\uparrow J^\downarrow \mathcal{X} = J^\downarrow x.$$

801 This – in turn – implies

$$J^\uparrow J^\downarrow \mathcal{X} = [J^\downarrow \Delta_{reg.} J^\uparrow - z Id]^{-1} J^\downarrow x.$$

802 Using

$$P_0^{high} \mathcal{X} = \mathcal{X},$$

803 we then have

$$\mathcal{X} = J^\uparrow [J^\downarrow \Delta_{reg.} J^\uparrow - z Id]^{-1} J^\downarrow x.$$

804 We have thus concluded the proof if we can prove that  $J^\downarrow \Delta_{reg.} J^\uparrow$  is the Laplacian corresponding to  
805 the graph  $\underline{G}$  defined in Definition B.1. But this is a straightforward calculation.  $\square$

806 As a corollary, we find

807 **Corollary B.4.** We have

$$R_z(\Delta)^k \rightarrow J^\uparrow R^k(\Delta) J^\downarrow$$

808 *Proof.* This follows directly from the fact that

$$J^\downarrow J^\uparrow = Id_G.$$

809

□

## 810 C Proof of Theorem 3.4

811 Here we prove Theorem 3.4, which we restate for convenience:

812 **Theorem C.1.** Fix  $\epsilon > 0$  and  $z < 0$ . For arbitrary functions  $g, h : [0, \infty] \rightarrow \mathbb{R}$  with  $\lim_{\lambda \rightarrow \infty} g(\lambda) =$   
 813 const. and  $\lim_{\lambda \rightarrow \infty} h(\lambda) = 0$ , there are filters  $f_{z,\theta}^0, f_{z,\theta}^I$  of Type-0 and Type-I respectively such that  
 814  $\|f_{z,\theta}^0 - g\|_\infty, \|f_{z,\theta}^I - h\|_\infty < \epsilon$ .

815 *Proof.* The Stone-Weierstrass theorem (see e.g. [40]) states that any sub-algebra of continuous  
 816 functions that are constant at infinity is already dense (in the topology of uniform convergence) if  
 817 this sub-algebra separates points.

818 Thus – using the Stone-Weierstrass Theorem – all we have to prove to establish the claim is that for  
 819 every pair of points  $x, y \geq 0$  there is a function  $f_\theta$  with

$$f_\theta(x) \neq f_\theta(y).$$

820 But this is clear since (for  $z < 0$ ) the function

$$\frac{1}{\cdot - z} : [0, \infty) \longrightarrow \mathbb{R}$$

821 (which generates the algebra of functions we consider) is already everywhere defined and injective.  
 822 □

## 823 D Stability Theory

824 Here we provide stability results to input- and edge-weight- perturbations for our architecture. For  
 825 convenience, we restate our layer-wise update rule here again:

826 Given a feature matrix  $X^\ell \in \mathbb{R}^{N \times F_\ell}$  in layer  $\ell$ , with column vectors  $\{X_j^\ell\}_{j=1}^{F_\ell}$ , the feature vector  
 827  $X_i^{\ell+1}$  in layer  $\ell+1$  is calculated as  $X_i^{\ell+1} = \text{ReLU}\left(\sum_{j=1}^{F_{\ell+1}} f_{z,\theta_{ij}^{\ell+1}}(\Delta) \cdot X_j^\ell + b_i^{\ell+1}\right)$  with a learnable  
 828 bias vector  $b_i^{\ell+1}$ . Collecting biases into a matrix  $B^{\ell+1} \in \mathbb{R}^{F_{\ell+1} \times N}$ , we efficiently implement this  
 829 using matrix-multiplications as

$$X^{\ell+1} = \text{ReLU}\left(\sum_{k=a}^K (T - \omega Id)^{-k} \cdot X^\ell \cdot W_k^{\ell+1} + B^{\ell+1}\right)$$

830 with weight matrices  $\{W_k^{\ell+1}\}$  in  $\mathbb{R}^{F_\ell \times F_{\ell+1}}$ . Biases are implemented as  $b_i = \beta_i \cdot \mathbf{1}_G$ , with  $\mathbf{1}_G$  the  
 831 vector of all ones on  $G$  and  $\beta_i \in \mathbb{R}$  learnable.

832 Our first result main-body of the paper then concerns stability to perturbations of input signals:

833 **Theorem D.1.** Let  $\Phi_L$  be the map associated to an  $L$ -layer deep ResolvNet. Denote the collection of  
 834 weight matrices in layer  $\ell$  by  $\mathscr{W}^\ell := \{W_k\}_{K=a}^{K_\ell}$ . We have

$$\|\Phi_L(X) - \Phi_L(Y)\|_2 \leq \|X - Y\|_2 \cdot \prod_{\ell=1}^L \|\mathscr{W}^\ell\|_z, \quad (23)$$

835 with

$$\|\mathscr{W}^\ell\|_z := \sum_{k=a}^K \frac{1}{|z|^k} \|W_k^\ell\|$$

836 aggregating singular values of weight matrices.

837 *Proof.* Let us denote (hidden) feature matrices in layer  $\ell$  by  $X^\ell$  (resp.  $Y^\ell$ ).

838 We note the following:

$$\begin{aligned}
\|X^L - Y^L\| &= \left\| \text{ReLU} \left( \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L + B^L \right) - \text{ReLU} \left( \sum_{k=a}^K R_z^k(\Delta) Y^{L-1} W_k^L + B^L \right) \right\| \\
&\leq \left\| \left( \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L + B^L \right) - \left( \sum_{k=a}^K R_z^k(\Delta) Y^{L-1} W_k^L + B^L \right) \right\| \\
&\leq \left\| \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L - \sum_{k=a}^K R_z^k(\Delta) Y^{L-1} W_k^L \right\| \\
&\leq \sum_{k=a}^K \|R_z^k(\Delta)\| \cdot \|X^{L-1} - Y^{L-1}\| \cdot \|W_k^L\| \\
&= \sum_{k=a}^K \frac{1}{|z|^k} \cdot \|X^{L-1} - Y^{L-1}\| \cdot \|W_k^L\| \\
&\leq \|\mathscr{W}^L\|_z \cdot \|X^{L-1} - Y^{L-1}\|.
\end{aligned}$$

839 Iterating through the layers yields the desired inequality (23).  $\square$

840 In preparation for our next result – Theorem D.5 below – we note the following:

841 **Lemma D.2.** Let  $\Phi_L$  be the map associated to an  $L$ -layer deep ResolvNet. With weights and biases  
842 denoted as above, we have

$$\|\Phi_L(X)\| \leq \|B^L\| + \sum_{m=0}^L \left( \prod_{j=0}^m \|\mathscr{W}^{L-1-j}\|_z \right) \|B^{L-1-j}\| + \left( \prod_{\ell=1}^L \|\mathscr{W}^\ell\|_z \right) \cdot \|X\|_2 \quad (24)$$

843 *Proof.* We have

$$\begin{aligned}
\|X\|^L &\leq \left\| \text{ReLU} \left( \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L + B^L \right) \right\| \\
&\leq \left\| \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L + B^L \right\| \\
&\leq \left\| \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L \right\| + \|B^L\| \\
&\leq \sum_{k=a}^K \|R_z^k(\Delta)\| \cdot \|X^{L-1}\| \cdot \|W_k^L\| + \|B^L\| \\
&\leq \left( \sum_{k=a}^K \frac{\|W_k^L\|}{|z|^k} \right) \cdot \|X^{L-1}\| + \|B^L\|.
\end{aligned}$$

844 Iterating this through all layers, we obtain (24).  $\square$

845 Before we can establish Theorem D.5 below, we need two additional (related) preliminary results:

846 **Lemma D.3.** Let us use the notation  $\tilde{R}_z := (\tilde{\Delta} - zId)^{-1}$  and  $R_z := (\Delta - zId)^{-1}$  for resolvents  
847 corresponding to two different Laplacians  $\Delta$  and  $\tilde{\Delta}$ . We have

$$\|R_z - \tilde{R}_z\| \leq \frac{1}{|z|^3} \|\Delta - \tilde{\Delta}\|$$

848 *Proof.* Let  $T$  and  $\tilde{T}$  be (finite dimensional) operators. Choose  $z$  so that it is neither an eigenvalue of  
849  $T$  nor  $\tilde{T}$ .

850 To showcase the principles underlying the proof, let us use the notation

$$R_z(T) \equiv \frac{1}{T-z}.$$

851 We note the following

$$\begin{aligned} & \frac{1}{\tilde{T}-z}(\tilde{T}-T)\frac{1}{T-z} \\ &= \frac{1}{\tilde{T}-z}\tilde{T}\frac{1}{T-z} - \frac{1}{\tilde{T}-z}T\frac{1}{T-z} \\ &= \left[ \frac{1}{\tilde{T}-z}(\tilde{T}-z) + \frac{z}{\tilde{T}-z} \right] \frac{1}{T-z} - \frac{1}{\tilde{T}-z} \left[ \frac{1}{T-z}(T-z) + \frac{z}{T-z} \right] \\ &= z \left( \frac{1}{T-z} - \frac{1}{\tilde{T}-z} \right). \end{aligned}$$

852 Rearranging and using

$$\|R_z(\Delta)\| = \|R_z(\tilde{\Delta})\| = \frac{1}{|z|}$$

853 together with the sub-multiplicativity of the operator-norm  $\|\cdot\|$  yields the claim.  $\square$

854 We also note the following estimate on differences of powers of resolvents:

855 **Lemma D.4.** Let  $\tilde{R}_z := (\tilde{\Delta} - zId)^{-1}$  and  $R_z := (\Delta - zId)^{-1}$ . For any natural number  $k$ , we have

$$\|\tilde{R}_z^k - R_z^k\| \leq \frac{k}{|z|^{k-1}} \|\tilde{R}_z - R_z\|$$

856 *Proof.* We note that for arbitrary matrices  $T, \tilde{T}$ , we have

$$\begin{aligned} \tilde{T}^k - T^k &= \tilde{T}^{k-1}(\tilde{T}-T) + (\tilde{T}^{k-1} - T^{k-1})T \\ &= \tilde{T}^{k-1}(\tilde{T}-T) + \tilde{T}^{k-2}(\tilde{T}-T)T + (\tilde{T}^{k-2} - T^{k-2})T^2. \end{aligned}$$

857 Iterating this and using

$$\|R_z(\Delta)\| = \|R_z(\tilde{\Delta})\| = \frac{1}{|z|}$$

858 for  $z < 0$  then yields the claim.  $\square$

859 Having established the preceding lemmata, we can now establish stability to perturbations of the edge  
860 weights:

861 **Theorem D.5.** Let  $\Phi_L$  and  $\tilde{\Phi}_L$  be the maps associated to ResolvNets with the same network archi-  
862 tecture, but based on Laplacians  $\Delta$  and  $\tilde{\Delta}$  respectively. We have

$$\|\Phi_L(X) - \tilde{\Phi}_L(X)\|_2 \leq (C_1(\mathcal{W}) \cdot \|X\|_2 + C_2(\mathcal{W}, \mathcal{B})) \cdot \|\Delta - \tilde{\Delta}\|. \quad (25)$$

863 Here, the stability constants  $C_1(\mathcal{W})$  and  $C_2(\mathcal{W}, \mathcal{B})$  are polynomials in (the largest) singular values  
864 of weight matrices and weight matrices as well as bias matrices, respectively.

865 *Proof.* Denote by  $X^\ell$  and  $\tilde{X}^\ell$  the (hidden) feature matrices generated in layer  $\ell$  for networks based  
866 on Laplacians  $\Delta$  and  $\tilde{\Delta}$  respectively: I.e. we have

$$X^\ell = \text{ReLU} \left( \sum_{k=a}^K R_z^k(\Delta) X^{\ell-1} W_k + B^\ell \right)$$

867 and

$$\tilde{X}^\ell = \text{ReLU} \left( \sum_{k=a}^K R_z^k(\tilde{\Delta}) \tilde{X}^{\ell-1} W_k + B^\ell \right).$$

868 Using the fact that  $\text{ReLU}(\cdot)$  is Lipschitz continuous with Lipschitz constant  $D = 1$ , we have

$$\begin{aligned} & \|X^L - \tilde{X}^L\| \\ &= \left\| \text{ReLU} \left( \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L + B^L \right) - \text{ReLU} \left( \sum_{k=a}^K R_z^k(\tilde{\Delta}) \tilde{X}^{L-1} W_k^L + B^L \right) \right\| \\ &\leq \left\| \left( \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L + B^L \right) - \left( \sum_{k=a}^K R_z^k(\tilde{\Delta}) \tilde{X}^{L-1} W_k^L + B^L \right) \right\| \\ &\leq \left\| \sum_{k=a}^K R_z^k(\Delta) X^{L-1} W_k^L - \sum_{k=a}^K R_z^k(\tilde{\Delta}) \tilde{X}^{L-1} W_k^L \right\| \\ &\leq \left\| \sum_{k=a}^K (R_z^k(\Delta) - R_z^k(\tilde{\Delta})) X^{L-1} W_k^L \right\| + \sum_{k=a}^K \|R_z(\tilde{\Delta})\| \cdot \|\tilde{X}^{L-1} - X^{L-1}\| \cdot \|W_k^L\| \\ &\leq \left\| \sum_{k=a}^K (R_z^k(\Delta) - R_z^k(\tilde{\Delta})) X^{L-1} W_k^L \right\| + \|\mathcal{W}^L\|_z \cdot \|\tilde{X}^{L-1} - X^{L-1}\| \\ &\leq \sum_{k=a}^K \left\| R_z^k(\Delta) - R_z^k(\tilde{\Delta}) \right\| \cdot \|X^{L-1}\| \cdot \|W_k^L\| + \|\mathcal{W}^L\|_z \cdot \|\tilde{X}^{L-1} - X^{L-1}\| \end{aligned}$$

869 Applying Lemma D.4 yields

$$\begin{aligned} & \|X^L - \tilde{X}^L\| \\ &\leq \left( \sum_{k=a}^K \frac{k}{|z|^{k-1}} \|W_k^L\| \right) \cdot \|X^{L-1}\| \cdot \|R_z(\Delta) - R_z(\tilde{\Delta})\| + \|\mathcal{W}^L\|_z \cdot \|\tilde{X}^{L-1} - X^{L-1}\|. \end{aligned}$$

870 Using Lemma D.3, we then have

$$\begin{aligned} & \|X^L - \tilde{X}^L\| \\ &\leq \left( \sum_{k=a}^K \frac{k}{|z|^{k+2}} \|W_k^L\| \right) \cdot \|X^{L-1}\| \cdot \|\Delta - \tilde{\Delta}\| + \|\mathcal{W}^L\|_z \cdot \|\tilde{X}^{L-1} - X^{L-1}\|. \end{aligned}$$

871 Lemma D.2 then yields

$$\begin{aligned} & \|X^L - \tilde{X}^L\| \\ &\leq \left( \sum_{k=a}^K \frac{k}{|z|^{k+2}} \|W_k^L\| \right) \cdot \left[ \|B^L\| + \sum_{m=0}^L \left( \prod_{j=0}^m \|\mathcal{W}^{L-1-k}\|_z \right) \|B^{L-1-k}\| + \left( \prod_{\ell=1}^L \|\mathcal{W}^\ell\|_z \right) \cdot \|X\|_2 \right] \cdot \|\tilde{\Delta} - \Delta\| \\ &+ \|\mathcal{W}^L\|_z \cdot \|\tilde{X}^{L-1} - X^{L-1}\|. \end{aligned}$$

872 Iterating this through the layers and collecting summands yields the desired relation (25).  $\square$

873 **E Stability under Scale Variations**

874 Here we provide details on the scale-invariance results discussed in Section 4.

875 In preparation, we will first need to prove a lemma relating powers of resolvents on the original graph  
876  $G$  and its limit-description  $\underline{G}$ :

877 **Lemma E.1.** Let  $\underline{R}_z := (\underline{\Delta} - zId)^{-1}$  and  $R_z := (\Delta - zId)^{-1}$ . For any natural number  $k$ , we have

$$\|J^\uparrow \underline{R}_z^k J^\downarrow - R_z^k\| \leq \frac{k}{|z|^{k-1}} \|J^\uparrow \underline{R}_z J^\downarrow - R_z\|$$

878 The proof proceeds in analogy to that of Lemma D.4:

879 *Proof.* We note that for arbitrary matrices  $T, \tilde{T}$ , we have

$$\begin{aligned} \tilde{T}^k - T^k &= \tilde{T}^{k-1}(\tilde{T} - T) + (\tilde{T}^{k-1} - T^{k-1})T \\ &= \tilde{T}^{k-1}(\tilde{T} - T) + \tilde{T}^{k-2}(\tilde{T} - T)T + (\tilde{T}^{k-2} - T^{k-2})T^2. \end{aligned}$$

880 Iterating this, using

$$\|R_z(\Delta)\| = \|J^\uparrow R_z(\underline{\Delta})J^\downarrow\| = \frac{1}{|z|}$$

881 for  $z < 0$  together with  $\|J^\uparrow\|, \|J^\downarrow\| \leq 1$  and

$$J^\uparrow \underline{R}_z^k J^\downarrow = (J^\uparrow \underline{R}_z J^\downarrow)^k$$

882 (which holds since  $J^\downarrow J^\uparrow = Id_{\underline{G}}$ ) then yields the claim.

883 Note that the equation

$$\|J^\uparrow R_z(\underline{\Delta})J^\downarrow\| = \frac{1}{|z|}$$

884 holds, because we may write

$$\|J^\uparrow R_z(\underline{\Delta})J^\downarrow\| = \left\| \lim_{\lambda_1(\Delta_{\text{high}}) \rightarrow \infty} R_z(\Delta) \right\| = \lim_{\lambda_1(\Delta_{\text{high}}) \rightarrow \infty} \|R_z(\Delta)\| = \lim_{\lambda_1(\Delta_{\text{high}}) \rightarrow \infty} \frac{1}{|z|} = \frac{1}{|z|}.$$

885 □

886 Hence let us now prove Stability-Theorem 4.1, which we restate here for convenience:

887 **Theorem E.2.** Let  $\Phi_L$  and  $\underline{\Phi}_L$  be the maps associated to ResolvNets with the same learned weight  
888 matrices and biases but deployed on graphs  $G$  and  $\underline{G}$  as defined in Section 2.2.2 . We have

$$\|\Phi_L(J^\uparrow \underline{X}) - J^\uparrow \Phi_L(\underline{X})\|_2 \leq (C_1(\mathcal{W}) \cdot \|X\|_2 + C_2(\mathcal{W}, \mathcal{B})) \cdot \|R_z(\Delta) - J^\uparrow R_z(\underline{\Delta})J^\downarrow\| \quad (26)$$

889 if the network is based on Type-0 resolvent filters (c.f. Section 3). Additionally, we have

$$\|\Phi_L(X) - J^\uparrow \Phi_L(J^\downarrow X)\|_2 \leq (C_1(\mathcal{W}) \cdot \|X\|_2 + C_2(\mathcal{W}, \mathcal{B})) \cdot \|R_z(\Delta) - J^\uparrow R_z(\underline{\Delta})J^\downarrow\| \quad (27)$$

890 if only Type-I filters are used in the network. Here  $C_1(\mathcal{W})$  and  $C_2(\mathcal{W}, \mathcal{B})$  are constants that depend  
891 polynomially on singular values of learned weight matrices  $\mathcal{W}$  and biases  $\mathcal{B}$ .

892 *Proof.* Let us first prove (27). To this end, let us define

$$\underline{X} := J^\downarrow X.$$

893 Let us further use the notation  $\underline{R}_z := (\underline{\Delta} - zId)^{-1}$  and  $R_z := (\Delta - zId)^{-1}$ .

894 Denote by  $X^\ell$  and  $\tilde{X}^\ell$  the (hidden) feature matrices generated in layer  $\ell$  for networks based on  
895 resolvents  $R_z$  and  $\underline{R}_z$  respectively: I.e. we have

$$X^\ell = \text{ReLU} \left( \sum_{k=a}^K R_z^k X^{\ell-1} W_k + B^\ell \right)$$

896 and

$$\tilde{X}^\ell = \text{ReLU} \left( \sum_{k=a}^K \underline{R}_z^k \tilde{X}^{\ell-1} W_k + \underline{B}^\ell \right).$$

897 Here, since bias terms are proportional to constant vectors on the graphs, as detailed in Section 3, we  
898 have

$$J^\downarrow B = \underline{B}$$

899 and

$$J^\uparrow \underline{B} = B \tag{28}$$

900 for bias matrices  $B$  and  $\underline{B}$  in networks deployed on  $G$  and  $\underline{G}$  respectively.

901 We then have

$$\begin{aligned} & \|\Phi_L(X) - J^\uparrow \Phi_L(J^\downarrow X)\| \\ &= \|X^L - J^\uparrow \tilde{X}^L\| \\ &= \left\| \text{ReLU} \left( \sum_{k=a}^K R_z^k X^{L-1} W_k^L + B^L \right) - J^\uparrow \text{ReLU} \left( \sum_{k=a}^K \underline{R}_z^k \tilde{X}^{L-1} W_k^L + \underline{B}^L \right) \right\| \\ &= \left\| \text{ReLU} \left( \sum_{k=a}^K R_z^k X^{L-1} W_k^L + B^L \right) - \text{ReLU} \left( \sum_{k=a}^K J^\uparrow \underline{R}_z^k \tilde{X}^{L-1} W_k^L + B^L \right) \right\|. \end{aligned}$$

902 Here we used the fact that since  $\text{ReLU}(\cdot)$  maps positive entries to positive entries and acts pointwise,  
903 it commutes with  $J^\uparrow$ . We also made use of (28).

904 Using the fact that  $\text{ReLU}(\cdot)$  is Lipschitz-continuous with Lipschitz constant  $D = 1$ , we can establish

$$\|\Phi_L(X) - J^\uparrow \Phi_L(J^\downarrow X)\| \leq \left\| \sum_{k=a}^K R_z^k X^{L-1} W_k^L - \sum_{k=a}^K J^\uparrow \underline{R}_z^k \tilde{X}^{L-1} W_k^L \right\|.$$

905 Using the fact that  $J^\downarrow J^\uparrow = Id_{\underline{G}}$ , we have

$$\|\Phi_L(X) - J^\uparrow \Phi_L(J^\downarrow X)\| \leq \left\| \sum_{k=1}^K R_z^k X^{L-1} W_k^L - \sum_{k=1}^K (J^\uparrow \underline{R}_z^k J^\downarrow) J^\uparrow \tilde{X}^{L-1} W_k^L \right\|.$$

906 From this, we find (using  $\|J^\uparrow\|, \|J^\downarrow\| \leq 1$ ), that

$$\begin{aligned} & \|X^L - J^\uparrow \tilde{X}^L\| \\ &\leq \left\| \sum_{k=0}^K R_z^k X^{L-1} W_k^L - \sum_{k=1}^K (J^\uparrow \underline{R}_z^k J^\downarrow) J^\uparrow \tilde{X}^{L-1} W_k^L \right\| \\ &\leq \left\| \sum_{k=1}^K (R_z^k - (J^\uparrow \underline{R}_z^k J^\downarrow)) X^{L-1} W_k^L \right\| + \sum_{k=1}^K \|J^\uparrow \underline{R}_z^k J^\downarrow\| \cdot \|J^\uparrow \tilde{X}^{L-1} - X^{L-1}\| \cdot \|W_k^L\| \\ &\leq \left\| \sum_{k=1}^K (R_z^k - (J^\uparrow \underline{R}_z^k J^\downarrow)) X^{L-1} W_k^L \right\| + \|\mathscr{W}^L\|_z \cdot \|J^\uparrow \tilde{X}^{L-1} - X^{L-1}\| \\ &\leq \sum_{k=1}^K \|R_z^k - (J^\uparrow \underline{R}_z^k J^\downarrow)\| \cdot \|X^{L-1}\| \cdot \|W_k^L\| + \|\mathscr{W}^L\|_z \cdot \|J^\uparrow \tilde{X}^{L-1} - X^{L-1}\| \end{aligned}$$

907 Applying Lemma E.1 yields

$$\begin{aligned} & \|X^L - J^\uparrow \tilde{X}^L\| \\ &\leq \left( \sum_{k=1}^K \frac{k}{|z|^{k-1}} \|W_k^L\| \right) \cdot \|R_z - (J^\uparrow \underline{R}_z J^\downarrow)\| \cdot \|X^{L-1}\| + \|\mathscr{W}^L\|_z \cdot \|J^\uparrow \tilde{X}^{L-1} - X^{L-1}\|. \end{aligned}$$



908 Lemma then D.2 in Appendix D established that we have

$$\|X^L\| \leq \|B^L\| + \sum_{m=0}^L \left( \prod_{j=0}^m \|\mathcal{W}^{L-1-k}\|_z \right) \|B^{L-1-k}\| + \left( \prod_{\ell=1}^L \|\mathcal{W}^\ell\|_z \right) \cdot \|X\|. \quad (29)$$

909 Hence the summand on the left-hand-side can be bounded in terms of a polynomial in singular values  
 910 of bias- and weight matrices, as well as  $\|X\|$  and most importantly the factor  $\|R_z - (J^\uparrow \underline{R}_z J^\downarrow)\|$   
 911 which tends to zero.

912 For the summand on the right-hand-side, we can iterate the above procedure (aggregating terms like  
 913 (29) multiplied by  $\|R_z - (J^\uparrow \underline{R}_z J^\downarrow)\|$ ) until reaching the last layer  $L = 1$ . There we observe

$$\begin{aligned} & \|X^1 - J^\uparrow \tilde{X}^1\| \\ &= \left\| \text{ReLU} \left( \sum_{k=1}^K R_z^k X W_k^1 + B^1 \right) - J^\uparrow \text{ReLU} \left( \sum_{k=1}^K \underline{R}_z^k J^\downarrow X W_k^1 + \underline{B}^1 \right) \right\| \\ &\leq \left\| \sum_{k=1}^K R_z^k X W_k^1 - \sum_{k=1}^K J^\uparrow \underline{R}_z^k J^\downarrow X W_k^1 \right\| \\ &\leq \left\| \sum_{k=1}^K (R_z^k - J^\uparrow \underline{R}_z^k J^\downarrow) X W_k^1 \right\| \\ &\leq \left( \sum_{k=1}^K \frac{k}{|z|^{k-1}} \|W_k^1\| \right) \cdot \|R_z - (J^\uparrow \underline{R}_z J^\downarrow)\| \cdot \|X\| \end{aligned}$$

914 The last step is only possible because we let the sums over powers of resolvents start at  $a = 1$  as  
 915 opposed to  $a = 0$ . In the latter case, there would have remained a term  $\|X - J^\uparrow J^\downarrow X\|$ , which would  
 916 not decay as  $\lambda_1(\Delta_{high}) \rightarrow \infty$ .

917 Aggregating terms, we build up the polynomial stability constants of (27) layer by layer, and  
 918 complete the proof.

919

920

921 The proof of (26) proceeds in complete analogy upon defining

$$X := J^\uparrow \underline{X}.$$

922 Note that starting with  $\underline{X}$  on  $\underline{G}$ , implies that we have

$$J^\uparrow J^\downarrow X \equiv J^\uparrow J^\downarrow (J^\uparrow \underline{X}) = J^\uparrow \underline{X} \equiv X.$$

923 This avoids any complications arising from employing Type-0 filters in this setting.

924

□

925 Next we transfer the previous result to the graph level setting:

926 **Theorem E.3.** Denote by  $\Psi$  the aggregation method introduced in Section 3. With  $\mu(G) = \sum_{i=1}^N \mu_i$   
 927 the total weight of the graph  $G$ , we have in the setting of Theorem 4.1 with Type-I filters, that

$$\|\Psi(\Phi_L(X)) - \Psi(\Phi_L(J^\downarrow X))\|_2 \leq \sqrt{\mu(G)} \cdot (C_1(\mathcal{W}) \cdot \|X\|_2 + C_2(\mathcal{W}, \mathcal{B})) \cdot \|R_z(\Delta) - J^\uparrow R_z(\underline{\Delta}) J^\downarrow\|.$$

928 *Proof.* Let us first recall that our aggregation scheme  $\Psi$  mapped a feature matrix  $X \in \mathbb{R}^{N \times F}$  to a  
 929 graph-level feature vector  $\Psi(X) \in \mathbb{R}^F$  defined component-wise as

$$\Psi(X)_j = \sum_{i=1}^N |X_{ij}| \cdot \mu_i.$$

930 In light of Theorem E.2, we are done with the proof, once we have established that

$$\|\Psi(\Phi_L(X)) - \Psi(\Phi_L(J^\downarrow X))\|_2 \leq \sqrt{\mu(G)} \cdot \|\Phi_L(X) - J^\uparrow \Phi_L(J^\downarrow X)\|_2.$$

931 To this end, we first note that

$$\Psi(J^\uparrow \underline{X}) = \Psi(\underline{X}).$$

932 Indeed, this follows from the fact that given a connected component  $R$  in  $G_{\text{high}}$ , the map  $J^\uparrow$  assigns  
 933 the same feature vector to each node  $r \in R \subseteq G$  (c.f. (20)), together with the fact that

$$\underline{\mu}_R = \sum_{r \in R} \mu_r.$$

934 Thus we have

$$\|\Psi(\Phi_L(X)) - \Psi(\underline{\Phi}_L(J^\downarrow X))\|_2 = \|\Psi(\Phi_L(X)) - \Psi(J^\uparrow \underline{\Phi}_L(J^\downarrow X))\|_2.$$

935 Next let us simplify notation and write

$$A = \Phi_L(X)$$

936 and

$$B = J^\uparrow \underline{\Phi}_L(J^\downarrow X)$$

937 with  $A, B \in \mathbb{R}^{N \times F}$ . We note:

$$\|\Psi(\Phi_L(X)) - \Psi(J^\uparrow \underline{\Phi}_L(J^\downarrow X))\|_2^2 = \sum_{j=1}^F \left( \sum_{i=1}^N (|A_{ij}| - |B_{ij}|) \cdot \mu_i \right)^2.$$

938 By means of the Cauchy-Schwarz inequality together with the inverse triangle-inequality, we have

$$\begin{aligned} \sum_{j=1}^F \left( \sum_{i=1}^N (|A_{ij}| - |B_{ij}|) \cdot \mu_i \right)^2 &\leq \sum_{j=1}^F \left[ \left( \sum_{i=1}^N |A_{ij} - B_{ij}|^2 \cdot \mu_i \right) \cdot \left( \sum_{i=1}^N \mu_i \right) \right] \\ &= \sum_{j=1}^F \left( \sum_{i=1}^N |A_{ij} - B_{ij}|^2 \cdot \mu_i \right) \cdot \mu(G). \end{aligned}$$

939 Since we have

$$\|\Phi_L(X) - J^\uparrow \underline{\Phi}_L(J^\downarrow X)\|_2^2 = \sum_{j=1}^F \left( \sum_{i=1}^N |A_{ij} - B_{ij}|^2 \cdot \mu_i \right),$$

940 the claim is established. □

## 941 **F Additional Details on Experiments:**

942 All experiments were performed on a single NVIDIA Quadro RTX 8000 graphics card.

### 943 **F.1 Node Classification**

944 **Datasets:** We test our approach for the task of node-classification on eight different standard  
 945 datasets across the entire homophily-spectrum. Among these, CITESEER [36], CORA-ML [25]  
 946 and PUBMED [26] are citation graphs. Here each node represents a paper and edges correspond  
 947 to citations. We also test on the MICROSOFT ACADEMIC graph [37] where an edge that is present  
 948 corresponds to co-authorship. Bag-of-word representations act as node features. The WEBKB  
 949 datasets CORNELL and TEXAS are datasets modeling links between websites at computer science  
 950 departments of various universities[29]. Node features are bag-of-words representation of the  
 951 respective web pages. We also consider the actor co-occurrence dataset ACTOR [39] as well as the  
 952 Wikipedia based dataset SQUIRREL [33].

953 **Experimental setup** We closely follow the experimental setup of [11] on which our codebase  
 954 builds: All models are trained for a fixed maximum (and unreachably high) number of  $n = 10000$   
 955 epochs. Early stopping is performed when the validation performance has not improved for 100  
 956 epochs. Test-results for the parameter set achieving the highest validation-accuracy are then reported.  
 957 Ties are broken by selecting the lowest loss (c.f. [42, 12]). Confidence intervals are calculated over  
 958 multiple splits and random seeds at the 95% confidence level via bootstrapping.

959 **Additional details on training and models:** We train all models on a fixed learning rate of

$$\text{lr} = 0.1.$$

960 Global dropout probability  $p$  of all models is optimized individually over

$$p \in \{0.3, 0.35, 0.4, 0.45, 0.5\}.$$

961 We use  $\ell^2$  weight decay and optimize the weight decay parameter  $\lambda$  for all models over

$$\lambda \in \{0.0001, 0.0005\}.$$

962 Where applicable (i.e. not for [12, 15]) we choose a two-layer deep convolutional architecture with  
963 the dimensions of hidden features optimized over

$$K_\ell \in \{32, 64, 128\}. \quad (30)$$

964 In addition to the hyperparameters specified above, some baselines have additional hyperparameters,  
965 which we detail here: BernNet uses an additional in-layer dropout rate of  $\text{dp\_rate} = 0.5$  and for its  
966 filters a polynomial order of  $K = 10$  as suggested in [15]. As suggested in [12], the hyperparameter  
967  $\alpha$  of PPNP is set to  $\alpha = 0.2$  on the MS\_ACADEMIC dataset and to  $\alpha = 0.1$  on other datasets.  
968 Hyperparameters depth  $T$  and number of stacks  $K$  of the ARMA convolutional layer [3] are set to  
969  $T = 1$  and  $K = 2$ . ChebNet also uses  $K = 2$  to avoid the known over-fitting issue [19] for higher  
970 polynomial orders. For MagNet we use  $K = 1$  as suggested in [47] and choose the parameter  $q$  as  
971 given in Table 1 of [47] for the respective datasets. The graph attention network [42] uses 8 attention  
972 heads, as suggested in [42].

973 For our ResolvNet model, we choose a depth of  $L = 1$  with hidden feature dimension optimized over  
974 the values in (30) as for baselines. We empirically observed in the setting of *unweighted* graphs, that  
975 rescaling the Laplacian as

$$\Delta_{nf} := \frac{1}{c_{nf}} \Delta$$

976 with a normalizing factor  $c_{nf}$  before calculating the resolvent

$$R_z(\Delta_{nf}) := (\Delta_{nf} - z \cdot Id)^{-1} \quad (31)$$

977 on which we base our ResolvNet architectures improved performance.

978 For our ResolvNet architecture, we express this normalizing factor in terms of the largest singular  
979 value  $\|\Delta\|$  of the (non-normalized) graph Laplacian. It is then selected among

$$c_{nf}/\|\Delta\| \in \{0.001, 0.01, 0.1, 2\}.$$

980 The value  $z$  in (31) is selected among

$$(-z) \in \{0.14, 0.15, 0.2, 0.25\}.$$

981 We base our ResolvNet architecture on Type-0 filters and choose the maximum resolvent-exponent  
982  $K$  as  $K = 1$ .

## 983 E.2 Graph Regression

984 **Datasets:** The first dataset we consider is the **QM7** dataset, introduced in [4, 35]. This dataset  
985 contains descriptions of 7165 organic molecules, each with up to seven heavy atoms, with all non-  
986 hydrogen atoms being considered heavy. A molecule is represented by its Coulomb matrix  $C^{\text{Clmb}}$ ,  
987 whose off-diagonal elements

$$C_{ij}^{\text{Clmb}} = \frac{Z_i Z_j}{|R_i - R_j|}$$

988 correspond to the Coulomb-repulsion between atoms  $i$  and  $j$ . We discard diagonal entries of Coulomb  
989 matrices; which would encode a polynomial fit of atomic energies to nuclear charge [35].

990 For each atom in any given molecular graph, the individual Cartesian coordinates  $R_i$  and the atomic  
991 charge  $Z_i$  are also accessible individually. To each molecule an atomization energy - calculated via  
992 density functional theory - is associated. The objective is to predict this quantity. The performance

Table 4: Targets of QM9

Symbol	Property	Unit
$U_0$	Internal energy at $0K$	$eV$
$U$	Internal energy at $298.15K$	$eV$
$H$	Enthalpy at $298.15K$	$eV$
$G$	Free energy at $298.15K$	$eV$
$U_0^{\text{ATOM}}$	Atomization energy at $0K$	$eV$
$U^{\text{ATOM}}$	Atomization energy at $298.15K$	$eV$
$H^{\text{ATOM}}$	Atomization enthalpy at $298.15K$	$eV$
$G^{\text{ATOM}}$	Atomization free energy at $298.15K$	$eV$
$c_v$	Heat capacity at $298.15K$	$\frac{\text{cal}}{\text{mol}\cdot\text{K}}$
$\mu$	Dipole moment	$D$
$\alpha$	Isotropic polarizability	$\alpha_0^3$
$\epsilon_{\text{HOMO}}$	Highest occupied molecular orbital energy	$eV$
$\epsilon_{\text{LUMO}}$	Lowest unoccupied molecular orbital energy	$eV$
$\Delta\epsilon$	Gap between $\epsilon_{\text{HOMO}}$ and $\epsilon_{\text{LUMO}}$	$eV$
$\langle R^2 \rangle$	Electronic spatial extent	$\alpha_0^2$
ZPVE	Zero point vibrational energy	$eV$
A	Rotational constant	$GHz$
B	Rotational constant	$GHz$
C	Rotational constant	$GHz$

metric is mean absolute error. Numerically, atomization energies are negative numbers in the range  $-600$  to  $-2200$ . The associated unit is  $[kcal/mol]$ .

The second dataset we consider is the **QM9** dataset [32], which consists of roughly 130 000 molecules in equilibrium. Beyond atomization energy, there are in total 19 targets available on **QM9**. We provide a complete list of targets together with abbreviations in Table 4 below:

Molecules in QM9 are not directly encoded via their Coulomb-matrices, as in QM7. However, positions and charges of individual molecules are available, from which the Coulomb matrix description is calculated for each molecule.

**Experimental Setup:** On both datasets, we randomly select 1500 molecules for testing and train on the remaining graphs. On QM7 we run experiments for 23 different random seeds and report mean and standard deviation. Due to computational limitations we run experiments for 3 different random seeds on the larger QM9 dataset, and report mean and standard deviation.

**Additional details on training and models:** All considered convolutional layers are incorporated into a two layer deep and fully connected graph convolutional architecture. In each hidden layer, we set the width (i.e. the hidden feature dimension) to

$$F_1 = F_2 = 64.$$

For BernNet, we set the polynomial order to  $K = 3$  to combat appearing numerical instabilities. ARMA is set to  $K = 2$  and  $T = 1$ . ChebNet uses  $K = 2$ . For all baselines, the standard mean-aggregation scheme is employed after the graph-convolutional layers to generate graph level features. Finally, predictions are generated via an MLP.

For our model, we choose a two-layer deep instantiation of our ResolvNet architecture introduced in Section 3. We choose Type-I filters and set  $z = -1$ . Laplacians are *not* rescaled and resolvents are thus given as

$$R_{-1}(\Delta) = (\Delta + Id)^{-1}.$$

As aggregation, we employ the graph level feature aggregation scheme introduced at the end of Section 3 with node weights set to atomic charges of individual atoms. Predictions are then generated via a final MLP with the same specifications as the one used for baselines.

All models are trained independently on each respective target.

1019 **Results:** Beyond the results already showcased in the main body of the paper, we here provide  
 1020 results for ResolvNet as well as baselines on all targets of Table 4. These results are collected in  
 1021 Table 5, Table 6 and Table 7 below.

1022 As is evident from the tables, the ResolvNet architecture produces mean-absolute-errors comparable  
 1023 to those of baselines on 1/4 of targets, while it performs significantly better on 3/4 of targets.

1024 The difference in performance is especially significant on the (extensive) energy targets of Table 5. In  
 1025 this Table, baselines are out-performed by factors varying between 4 and 15.

1026 Table 6 contains three additional targets where MAEs produced by ResolvNet are lower by factors  
 1027 varying between roughly two and four, when compared to baselines.

1028 Table 7 finally contains MAEs corresponding to predictions of rotational constants. Here our model  
 1029 yields a comparable error on one target and provides better results than baselines on two out of three  
 1030 targets.

Table 5: Energy prediction MAEs [ $eV$ ]. Our Model is marked **R.N.** for **ResolvNet**.

Property	$U_0$	$U$	$H$	$G$	$U_0^{\text{ATOM}}$	$U^{\text{ATOM}}$	$H^{\text{ATOM}}$	$G^{\text{ATOM}}$
BernNet	370.42±38.91	382.64±36.52	398.32±46.00	362.69±24.84	3.112±0.285	3.096±0.249	3.046±0.277	2.919 ±0.375
GCN	381.41±0.42	376.41±7.10	368.01±16.77	380.65±6.67	2.766±0.081	2.828±0.091	2.803±0.077	2.575±0.084
ChebNet	345.74±12.30	346.39±19.11	398.32±22.48	350.22±12.32	2.665±0.040	2.672±0.056	2.745±0.104	2.477±0.036
ARMA	327.62±19.83	316.09±18.06	322.74±16.32	320.72±11.98	2.588±0.117	2.570±0.088	2.600±0.096	2.326±0.101
R.N.	<b>21.72</b> ±5.79	<b>19.14</b> ±7.19	<b>31.18</b> ±8.622	<b>53.50</b> ±4.58	<b>0.605</b> ±0.015	<b>0.588</b> ±0.024	<b>0.593</b> ±0.025	<b>0.607</b> ±0.041

Table 6: Various target prediction MAEs. Our Model is marked **R.N.** for **ResolvNet**.

Property	$c_v$ [ $\frac{\text{cal}}{\text{mol}\cdot\text{K}}$ ]	$\mu$ [ $D$ ]	$\alpha$ [ $\alpha_0^3$ ]	$\epsilon_{\text{HOMO}}$ [ $eV$ ]	$\epsilon_{\text{LUMO}}$ [ $eV$ ]	$\Delta\epsilon$ [ $eV$ ]	$\langle R^2 \rangle$ [ $\alpha_0^2$ ]	ZPVE [ $eV$ ]
BernNet	2.610±0.986	0.948±0.042	3.519±0.288	0.376±0.028	0.649±0.092	0.841±0.085	157.982 ±34.804	0.237 ±0.032
GCN	1.521±0.038	0.936±0.003	3.114±0.112	0.301±0.009	0.523±0.018	0.566±0.016	130.461±5.445	0.185±0.004
ChebNet	1.455±0.053	0.881±0.007	3.049±0.092	0.234±0.005	0.433±0.018	0.515±0.010	132.695±2.218	0.180±0.005
ARMA	1.327±0.034	0.806±0.031	2.676±0.087	<b>0.228</b> ±0.010	<b>0.333</b> ±0.009	<b>0.380</b> ±0.007	<b>93.760</b> ±4.122	0.152±0.006
R.N.	<b>0.747</b> ±0.015	<b>0.776</b> ±0.018	<b>1.308</b> ±0.034	0.313±0.002	0.423±0.011	0.531±0.016	97.614±2.308	<b>0.041</b> ±0.008

Table 7: Rotational constants prediction MAEs. Our Model is marked **R.N.** for **ResolvNet**.

Property	$A$ [ $GHz$ ]	$B$ [ $GHz$ ]	$C$ [ $GHz$ ]
BernNet	0.888±0.034	0.342±0.002	0.243±0.002
GCN	0.848±0.027	0.281±0.004	0.183±0.002
ChebNet	0.797±0.034	0.262±0.003	0.171±0.003
ARMA	<b>0.715</b> ±0.017	0.259±0.004	0.168±0.004
R.N.	0.783±0.802	<b>0.249</b> ±0.002	<b>0.158</b> ±0.001

### 1031 F.3 Scale Invariance

1032 **Dataset:** Again, we make use of the QM7 dataset [35] and its Coulomb matrix description

$$C_{ij}^{\text{Cmb}} = \frac{Z_i Z_j}{|R_i - R_j|} \quad (32)$$

1033 of molecules.

1034 **Details on collapsing procedure:** We modify (all) molecular graphs in QM7 by deflecting hydro-  
 1035 gen atoms (H) out of their equilibrium positions towards the respective nearest heavy atom. This is  
 1036 possible since the QM7 dataset also contains the Cartesian coordinates of individual atoms.

1037 This introduces a two-scale setting precisely as discussed in section 2: Edge weights between heavy  
 1038 atoms remain the same, while Coulomb repulsions between H-atoms and respective nearest heavy  
 1039 atom increasingly diverge; as is evident from (32).

1040 Given an original molecular graph  $G$  with node weights  $\mu_i = Z_i$ , the corresponding limit graph  
 1041  $\underline{G}$  corresponds to a coarse grained description, where heavy atoms and surrounding H-atoms are  
 1042 aggregated into single super-nodes in the sense of Section 2.2.2 .

1043 Mathematically,  $\underline{G}$  is obtained by removing all nodes corresponding to H-atoms from  $G$ , while adding  
 1044 the corresponding charges  $Z_H = 1$  to the node-weights of the respective nearest heavy atom. Charges  
 1045 in (32) are modified similarly to generate the weight matrix  $\underline{W}$ .

1046 On original molecular graphs, atomic charges are provided via one-hot encodings. For the graph of  
 1047 methane – consisting of one carbon atom with charge  $Z_C = 6$  and four hydrogen atoms of charges  
 1048  $Z_H = 1$  – the corresponding node-feature-matrix is e.g. given as

$$X = \begin{pmatrix} 0 & 0 & \cdots & 0 & 1 & 0 \cdots \\ 1 & 0 & \cdots & 0 & 0 & 0 \cdots \\ 1 & 0 & \cdots & 0 & 0 & 0 \cdots \\ 1 & 0 & \cdots & 0 & 0 & 0 \cdots \\ 1 & 0 & \cdots & 0 & 0 & 0 \cdots \end{pmatrix}$$

1049 with the non-zero entry in the first row being in the 6<sup>th</sup> column, in order to encode the charge  $Z_C = 6$   
 1050 for carbon.

1051 The feature vector of an aggregated node represents charges of the heavy atom and its neighbouring  
 1052 H-atoms jointly.

1053 As discussed in Definition 3.2, node feature matrices are translated as  $\underline{X} = J^\downarrow X$ . Applying  $J^\downarrow$   
 1054 to one-hot encoded atomic charges yields (normalized) bag-of-word embeddings on  $\underline{G}$ : Individual  
 1055 entries of feature vectors encode how much of the total charge of the super-node is contributed by  
 1056 individual atom-types. In the example of methane, the limit graph  $\underline{G}$  consists of a single node with  
 1057 node-weight

$$\mu = 6 + 1 + 1 + 1 + 1 = 10.$$

1058 The feature matrix

$$\underline{X} = J^\downarrow X$$

1059 is a single row-vector given as

$$\underline{X} = \left( \frac{4}{10}, 0, \dots, 0, \frac{6}{10}, 0, \dots \right).$$

1060 **Results:**

For convenience, we repeat here in Table 8 and Figure 12 the results corresponding to the use of resolution-limited data in the form of coarse-grained molecular graphs during inference, that were already presented in the main body of the paper.

Table 8: MAE on QM7 via coarsified molecular graphs.

1061

QM7	MAE [kcal/mol]
BernNet	580.67 $\pm$ 99.27
GCN	124.53 $\pm$ 34.58
ChebNet	645.14 $\pm$ 34.59
ARMA	248.96 $\pm$ 15.56
ResolvNet	<b>16.23<math>\pm</math>2.74</b>

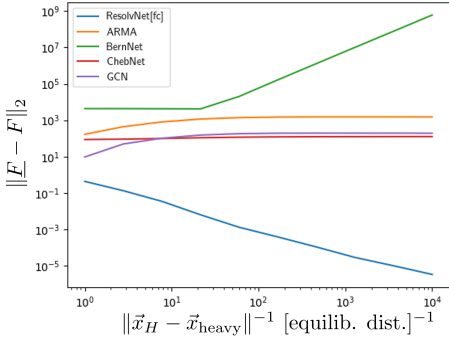


Figure 12: Feature-vector-difference for collapsed ( $\underline{F}$ ) and deformed ( $F$ ) graphs.

1062 **G Analysis of Computational Overhead**

1063 Here we provide an analysis of the overhead of our ResolvNet method. As is evident from Tables 9,  
 1064 10, 11 below, on most datasets our method is not the most memory intensive to train when compared  
 1065 to representative (spatial and spectral) baselines. For training times (total and per-epoch), we note  
 1066 that on most small to medium sized graphs, our model is not the slowest to train. On larger graphs it  
 1067 does take longer to train. Regarding complexity, the node update for our model is essentially  $\mathcal{O}(N^2)$   
 1068 (dense-dense matrix multiplication), while message passing baselines scale linearly in the number of  
 1069 edges.

Table 9: Maximal Memory Consumption [GB] while training a single model of depth 2 and width 32 for learning rate  $\text{lr} = 0.1$ , dropout  $p = 0.5$ , weight decay  $\lambda = 10^{-4}$  and early stopping patience  $t = 100$ . All measurements performed on the same GPU via `torch.cuda.max_memory_allocated()`.

	MS_Acad.	Cora	Pubmed	Citeseer	Cornell	Actor	Squirrel	Texas
ResolvNet	3.47	0.1266	2.9915	0.0996	0.0070	0.4936	0.2915	0.0175
GAT	1.49	0.1559	0.6486	0.1105	0.0228	0.3666	2.1107	0.0219
ChebNet	10.19	0.4741	0.4848	0.3389	0.0249	0.4830	6.3569	0.0241

1070  
1071  
1072

Table 10: Training Time [s] for training a single model of depth 2 and width 32 for learning rate  $\text{lr} = 0.1$ , dropout  $p = 0.5$ , weight decay  $\lambda = 10^{-4}$  and early stopping patience  $t = 100$ . All measurements performed on the same GPU.

	MS_Acad.	Cora	Pubmed	Citeseer	Cornell	Actor	Squirrel	Texas
ResolvNet	474.409	3.671	34.140	1.387	1.745	9.623	4.874	0.875
GAT	34.388	2.194	5.741	0.891	2.123	1.610	23.060	1.375
ChebNet	87.567	6.818	3.221	2.833	2.713	1.488	14.383	4.511

1073  
1074  
1075

Table 11: Average Training Time per Epoch [ms] for training a single model of depth 2 and width 32 for learning rate  $\text{lr} = 0.1$ , dropout  $p = 0.5$ , weight decay  $\lambda = 10^{-4}$  and early stopping patience  $t = 100$ . All measurements performed on the same GPU.

	MS_Acad.	Cora	Pubmed	Citeseer	Cornell	Actor	Squirrel	Texas
ResolvNet	1359.34	13.16	161.80	11.01	2.58	32.51	41.30	2.54
GAT	60.01	8.22	29.59	7.24	3.93	15.05	62.49	4.07
ChebNet	202.23	12.11	14.31	10.61	3.89	13.28	126.17	3.83