
HoloNets: Spectral Convolutions do extend to Directed Graphs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Within the graph learning community, conventional wisdom dictates that spectral
2 convolutional networks may only be deployed on undirected graphs: Only there
3 could the existence of a well-defined graph Fourier transform be guaranteed, so
4 that information may be translated between spatial- and spectral domains. Here
5 we show this traditional reliance on the graph Fourier transform to be superfluous
6 and – making use of certain advanced tools from complex analysis and spectral
7 theory – extend spectral convolutions to directed graphs. We provide a frequency-
8 response interpretation of newly developed filters, investigate the influence of the
9 basis used to express filters and discuss the interplay with characteristic operators
10 on which networks are based. In order to thoroughly test the developed theory,
11 we conduct experiments in real world settings, showcasing that directed spectral
12 convolutional networks provide new state of the art results for heterophilic node
13 classification on many datasets and – as opposed to baselines – may be rendered
14 stable to resolution-scale varying topological perturbations.

15 1 Introduction

16 A particularly prominent line of research for graph neural networks is that of spectral convolutional
17 architectures. These are among the theoretically best understood graph learning methods [34, 48, 32]
18 and continue to set the state of the art on a diverse selection of tasks [6, 23, 24, 57]. Furthermore,
19 spectral interpretations allow to better analyse expressivity [3], shed light on shortcomings of
20 established models [43] and guide the design of novel methods [8].

21 Traditionally, spectral convolutional filters are defined making use of the notion of a graph Fourier
22 transform: Fixing a self-adjoint operator on an undirected N -node graph – traditionally a suitably
23 normalized graph Laplacian $L = U^\top \Lambda U$ with eigenvalues $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$ – a notion of Fourier
24 transform is defined by projecting a given signal x onto the eigenvectors of L via $x \mapsto Ux$. Since L
25 is self-adjoint, the eigenvectors form a complete basis and no information is lost in the process.

26 In analogy with the Euclidean convolution theorem, early spectral networks then defined convolution
27 as multiplication in the "graph-Fourier domain" via $x \mapsto U^\top \cdot \text{diag}(\theta_1, \dots, \theta_N) \cdot Ux$, with learnable
28 parameters $\{\theta_1, \dots, \theta_N\}$ [9]. To avoid calculating an expensive explicit eigendecomposition U , [15]
29 proposed to instead parametrize graph convolutions via $x \mapsto U^\top g_\theta(\Lambda) Ux$, with g_θ a learnable
30 scalar function applied to the eigenvalues Λ as $g_\theta(\Lambda) = \text{diag}(g_\theta(\lambda_1), \dots, g_\theta(\lambda_N))$. This precisely
31 reproduces the mathematical definition of applying a scalar function g_θ to a self-adjoint operator L ,
32 so that choosing g_θ to be a (learnable) polynomial allowed to implement filters computationally much
33 more economically as $g_\theta(L) = \sum_{k=1}^K \theta_k L^k$. Follow up works then considered the influence of the
34 basis in which filters $\{g_\theta\}$ are learned [23, 35, 58] and established that such filters provide networks
35 with the ability to generalize to unseen graphs [34, 49, 32].

36 Common among all these works, is the need for the underlying graph to be undirected: Only then are
 37 the characteristic operators self-adjoint, so that a complete set of eigenvectors exists and the graph
 38 Fourier transform U – used to define the filter $g_\theta(L)$ via $x \mapsto Ug_\theta(\Lambda)U^\top x$ – is well-defined.

39 Currently however, the graph learning community is endeavouring to finally also account for the
 40 previously neglected directionality of edges, when designing new methods [59, 47, 5, 18, 25]. Since
 41 characteristic operators on digraphs are generically not self-adjoint, traditional spectral approaches so
 42 far remained inaccessible in this undertaking. Instead, works such as [59, 25] resorted to limiting
 43 themselves to certain specialized operators able to preserve self-adjointness in this directed setting.
 44 While this approach is not without merit, the traditional adherence to the graph Fourier transform
 45 remains a severely limiting factor when attempting to extend spectral networks to directed graphs.

46 **Contributions:** In this paper we argue to completely dispose with this reliance on the graph Fourier
 47 transform and instead take the concept of learnable functions applied to characteristic operators as
 48 fundamental. This conceptual shift allows us to consistently define spectral convolutional filters on
 49 directed graphs. We provide a corresponding frequency perspective, analyze the interplay with chosen
 50 characteristic operators and discuss the importance of the basis used to express these novel filters.
 51 The developed theory is thoroughly tested on real world data: It is found that that directed spectral
 52 convolutional networks provide new state of the art results for heterophilic node classification and – as
 53 opposed to baselines – may be rendered stable to resolution-scale varying topological perturbations.

54 2 Signal processing on directed Graphs

55 **Weighted directed graphs:** A directed graph $G := (\mathcal{G}, \mathcal{E})$ is a collection of nodes \mathcal{G} and edges
 56 $\mathcal{E} \subseteq \mathcal{G} \times \mathcal{G}$ for which $(i, j) \in \mathcal{E}$ does not necessarily imply $(j, i) \in \mathcal{E}$. We allow nodes $i \in \mathcal{G}$ to have
 57 individual node-weights $\mu_i > 0$ and generically assume edge-weights $w_{ij} \geq 0$ not necessarily equal
 58 to unity or zero. In a social network, a **node weight** $\mu_i = 1$ might signify that a node represents a
 59 single user, while a weight $\mu_j > 1$ would indicate that node j represents a group of users. Similarly,
 60 **edge weights** $\{w_{ij}\}$ could be used to encode how many messages have been exchanged between
 61 nodes i and j . Importantly, since we consider directed graphs, we generically have $w_{ij} \neq w_{ji}$.

Edge weights also determine the so called **reaches** of a graph, which generalize the concept of connected components of undirected graphs [54]: A subgraph $R \subseteq G$ is called reach, if for any two vertices $a, b \in R$ there is a directed path in R along which the (directed) edge weights do not vanish, and R simultaneously possesses no outgoing connections (i.e. for any $c \in G$ with $c \notin R$: $w_{ca} = 0$). For us, this concept will be important in generalizing the notion of scale insensitive networks [33] to directed graphs in Section 3.3 below.

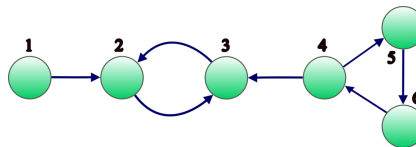


Figure 1: A di-graph with reaches $R_1 = \{1, 2, 3\}$ and $R_2 = \{3, 4, 5, 6, 2\}$.

63 **Feature spaces:** Given F -dimensional node features on a graph with $N = |\mathcal{G}|$ nodes, we may
 64 collect individual node-feature vectors into a feature matrix X of dimension $N \times F$. Taking into
 65 account our node weights, we equip the space of such signals with an inner-product according
 66 to $\langle X, Y \rangle = \text{Tr}(X^*MY) = \sum_{i=1}^N \sum_{j=1}^F (\bar{X}_{ij}Y_{ij})\mu_i$ with $M = \text{diag}(\{\mu_i\})$ the diagonal matrix
 67 of node-weights. Here X^* denotes the (hermitian) adjoint of X (c.f. Appendix B for a brief
 68 recapitulation). Associated to this inner product is the standard 2-norm $\|X\|_2^2 = \sum_{i=1}^N \sum_{j=1}^F |X_{ij}|^2 \mu_i$.

69 **Characteristic Operators:** Information about the geometry of a graph is encapsulated into the set
 70 of edge weights, collected into the weight matrix W . From this, the diagonal in-degree and out-degree
 71 matrices ($D_{ii}^{\text{in}} = \sum_j W_{ij}$, $D_{jj}^{\text{out}} = \sum_i W_{ij}$) may be derived. Together with the the node-weight matrix
 72 M defined above, various characteristic operators capturing the underlying geometry of the graph
 73 may then be constructed. Relevant to us – apart from the weight matrix W – will especially be the
 74 (in-degree) Laplacian $L^{\text{in}} := M^{-1}(D^{\text{in}} - W)$, which is intimately related to consensus and diffusion
 75 on directed graphs [55]. Importantly, such characteristic operators T are generically not self-adjoint.
 76 Hence they do not admit a complete set of eigenvectors and their spectrum $\sigma(T)$ contains complex
 77 eigenvalues $\lambda \in \mathbb{C}$. Appendix B contains additional details on such operators, their canonical (Jordan)
 78 decomposition and associated *generalized* eigenvectors.

79 **3 Spectral Convolutions on directed graphs**

80 Since characteristic operators on directed graphs generically do not admit a complete set of orthogonal
 81 eigenvectors, we cannot make use of the notion of a graph Fourier transform to consistently define
 82 filters of the form $g_\theta(T)$. While this might initially seem to constitute an insurmountable obstacle,
 83 the task of defining operators of the form $g(T)$ for a given operator T and appropriate classes of
 84 scalar-valued functions $\{g\}$ – such that relations between the functions $\{g\}$ translate into according
 85 relation of the operators $\{g(T)\}$ – is in fact a well studied problem [20, 13]. Corresponding techniques
 86 typically bear the name "functional calculus" and importantly are also definable if the underlying
 87 operator T is not self-adjoint [14]:

88 **3.1 The holomorphic functional calculus**

89 In the undirected setting, it was possible to essentially apply arbitrary functions $\{g\}$ to the characteris-
 90 tic operator $T = U^\top \Lambda U$ by making use of the complete eigendecomposition as $g(T) := U^\top g_\theta(\Lambda) U$.
 91 However, a different approach to consistently defining the matrix $g(T)$ – not contingent on such a
 92 decomposition – is available if one restricts g to be a **holomorphic** function: For a given subset U
 93 of the complex plane, these are the complex valued functions $g : U \rightarrow \mathbb{C}$ for which the complex
 94 derivative $dg(z)/dz$ exists everywhere on the domain U (c.f. Appendix D for more details).

The property of such holomorphic functions that is exploited in order to consistently define the matrix $g(T)$ is the fact that any function value $g(\lambda)$ can be reproduced by calculating an integral of the function g along a path Γ encircling λ (c.f. also Fig. 2) as

$$g(\lambda) = -\frac{1}{2\pi i} \oint_{\Gamma} g(z) \cdot (\lambda - z)^{-1} dz. \quad (1)$$

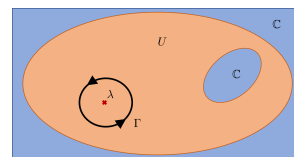


Figure 2: Cauchy Integral (1)

96 In order to define the matrix $g(T)$, the formal replacement $\lambda \mapsto T$ is then made on both sides of (1),
 97 with the path Γ now not only encircling a single value λ but all eigenvalues $\lambda \in \sigma(T)$ (c.f. also Fig. 3):

$$g(T) := -\frac{1}{2\pi i} \oint_{\Gamma} g(z) \cdot (T - z \cdot Id)^{-1} dz \quad (2)$$

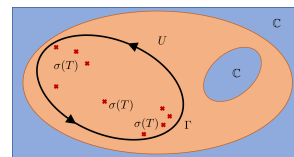


Figure 3: Operator Integral (2)

98 Note that $(T - z \cdot Id)^{-1}$ – and hence the integral in (2) – is indeed well-defined: All eigenvalues of T are assumed to lie *inside* the path Γ . For any choice of integration variable z *on* this path Γ , the matrix $(T - z \cdot Id)$ is thus indeed invertible, since z is never an eigenvalue.

99 The integral in (2) defines what is called the **holomorphic functional calculus** [19, 30]. Importantly,
 100 this holomorphic functional calculus (2) agrees with algebraic relations: Applying a polynomial
 101 $g(\lambda) := \sum_{k=0}^K a_k \lambda^k$ to T yields $g(T) = \sum_{k=0}^K a_k T^k$. Similarly applying the function $g(\lambda) = 1/\lambda$
 102 yields $g(T) = T^{-1}$ provided T is invertible. Appendix E details the calculations.

103 **3.2 Spectral convolutional filters on directed graphs**

104 Since the holomorphic functional calculus is evidently no longer contingent on T being self-adjoint,
 105 it indeed provides an avenue to consistently define spectral convolutional filters on directed graphs.

106 **Parametrized spectral convolutional filters:** In practice it is of course prohibitively expensive to
 107 continuously compute the integral (2) as the learnable function g is updated during training. Instead,
 108 we propose to represent a generic holomorphic function g via a set of basis functions $\{\Psi_i\}_{i \in I}$ as
 109 $g_\theta(z) := \sum_{i \in I} \theta_i \cdot \Psi_i(z)$ with learnable coefficients $\{\theta_i\}_{i \in I}$ parametrizing the filter g_θ .

110 For the 'simpler' basis functions $\{\Psi_i\}_{i \in I}$, we either precompute the integral (2), or perform it
 111 analytically (c.f. Section 3.3 below). During training and inference the matrices $\Psi_i(T) \equiv$
 112 $-\frac{1}{2\pi i} \oint_{\Gamma} \Psi_i(z) \cdot (T - z \cdot Id)^{-1} dz$ are then already computed and learnable filters are given as

$$g_\theta(T) := \sum_{i \in I} \theta_i \cdot \Psi_i(T).$$

114 Generically, each coefficient θ_i may be chosen as a *complex* number; equivalent to two *real* parameters.
 115 If the functions $\{\Psi_i\}_{i \in I}$ are chosen such that each matrix $\Psi_i(T)$ contains only real entries (e.g. for Ψ a
 116 polynomial with real coefficients), it is possible to restrict convolutional filters to being purely real: In
 117 this setting, choosing the parameters $\{\theta_i\}$ to be purely real as well, leads to $g_\theta(T) = \sum_{i \in I} \theta_i \cdot \Psi_i(T)$
 118 itself being a matrix that contains only real entries. In this way, complex numbers need never to
 119 appear within our network, if this is not desired. In Theorem 4.1 of Section 4 below, we discuss how,
 120 under mild and reasonable assumptions, such a complexity-reduction to using only real parameters
 121 can be performed without decreasing the expressive power of corresponding networks.

122 Irrespective of whether real or complex weights are employed, the utilized filter bank $\{\Psi_i\}_{i \in I}$
 123 determines the space of learnable functions $g_\theta \in \text{span}(\{\Psi_i\}_{i \in I})$ and thus contributes significantly to
 124 the inductive bias present in the network. It should thus be adjusted to the respective task at hand.

125 **The Action of Filters in the Spectral Domain:** In order to determine which basis functions are
 126 adapted to which tasks, a "frequency-response" interpretation of spectral filters is expedient:

127 In the **undirected setting** this proceeded by decomposing any characteristic operator T into a sum
 128 $T = \sum_{\lambda \in \sigma(T)} \lambda \cdot P_\lambda$ over its distinct eigenvalues. The spectral action of any function g was then
 129 given by $g(T) = \sum_{\lambda \in \sigma(T)} g(\lambda) \cdot P_\lambda$. Here the spectral projections P_λ project each vector to the space
 130 spanned by the eigenvectors $\{v_i\}$ corresponding to the eigenvalue λ (i.e. satisfying $(T - \lambda Id)v_i = 0$).

131 In the **directed setting**, there only exists a basis of *generalized* eigenvectors $\{w_i\}_{i=1}^N$; each satisfying
 132 $(T - \lambda \cdot Id)^m w_i = 0$ for some $\lambda \in \sigma(T)$ and $m \in \mathbb{N}$ (c.f. Appendix B). Denoting by P_λ the matrix
 133 projecting onto the space spanned by these *generalized* eigenvectors associated to the eigenvalue
 134 $\lambda \in \sigma(T)$, any operator T may be written as¹ $T = \sum_{\lambda \in \sigma(T)} \lambda \cdot P_\lambda + \sum_{\lambda \in \sigma(T)} (T - \lambda \cdot Id) \cdot P_\lambda$. It
 135 can then be shown [30], that the spectral action of a given function g is given as

$$g(T) = \sum_{\lambda \in \sigma(T)} g(\lambda) P_\lambda + \sum_{\lambda \in \sigma(T)} \left[\sum_{n=1}^{m_\lambda-1} \frac{g^{(n)}(\lambda)}{n!} (T - \lambda \cdot Id)^n \right] P_\lambda. \quad (3)$$

136 Here the number m_λ is the algebraic multiplicity of the eigenvalue λ ; i.e. the dimension of the
 137 associated *generalized* eigenspace. The notation $g^{(n)}$ denotes the n^{th} complex derivative of g . The
 138 appearance of such derivative terms in (3) is again evidence, that we indeed needed to restrict from
 139 generic- to differentiable functions² in order to sensibly define directed spectral convolutional filters.

140 It is instructive to gain some intuition about the second sum on the right-hand-side of the frequency
 141 response (3), as it is not familiar from undirected graphs (since it vanishes if T is self-adjoint):
 As an example consider the un-weighted directed path graph on three nodes
 depicted in Fig. 4 and choose as characteristic operator T the adjacency matrix
 (i.e. $T = W$). It is not hard to see (c.f. Appendix C for an explicit calculation)
 142 that the only eigenvalue of W is given by $\lambda = 0$ with algebraic multiplicity
 $m_\lambda = 3$. Since spectral projections always satisfy $\sum_{\lambda \in \sigma(T)} P_\lambda = Id$ (c.f.
 Appendix B), and here $\sigma(W) = \{0\}$ we thus have $P_{\lambda=0} = Id$ in this case.

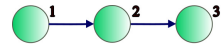


Figure 4: Directed path on 3 nodes

143 Suppose now we are tasked with finding a (non-trivial) holomorphic filter $g(\lambda)$ such that $g(T) = 0$.
 144 The right-hand sum in (3) implies, that beyond $g(0) = 0$, also the first and second derivative of $g(\lambda)$
 145 needs to vanish at $\lambda = 0$ to achieve this. Hence the zero of $g(\lambda)$ at $\lambda = 0$ must be at least of order
 146 three; or equivalently for $\lambda \rightarrow 0$ we need $g(\lambda) = o(\lambda^3)$. This behaviour is of course exactly mirrored
 147 in the spatial domain: As applying W simply moves information at a given node along the path,
 148 applying W once or twice still leaves information present. After two applications, only node 3 still
 149 contains information and thus applying W^k precisely removes all information if and only if $k \geq 3$.

150 Without the assumption of acyclicity, the spectrum of characteristic operators of course generically
 151 does not consist only of the eigenvalue $\lambda = 0$. Thus generically $P_{\lambda=0} \neq Id$ and the role played by
 152 the operator $T = W$ in the considerations above is instead played by its restriction $(T \cdot P_{\lambda=0})$ to the
 153 generalized eigenspace corresponding to the eigenvalue $\lambda = 0$.

154 For us, the spectral response (3) will provide guidance when designing and discussing scale-insensitive
 155 convolutional filters and corresponding networks on *directed* graphs in Sections 3.3 and 4 below.

¹Additional details on this so called Jordan Chevalley decomposition are provided in Appendix B.

²N.B.: A once-complex-differentiable function is automatically infinitely often differentiable [1].

156 **3.3 Explicit Filter Banks**

157 Having laid the theoretical foundations, we consider examples of task-adapted filter banks $\{\Psi_i\}_{i \in I}$.

158 **3.3.1 Bounded Spectral Domain: Faber Polynomials**

159 First, let us consider spectral networks on a single graph with a fixed characteristic operator T . From
 160 the holomorphic functional calculus (2), we infer that convolutional filters $\{g(T)\}$ are in principle
 161 provided by all holomorphic functions $\{g\}$ defined on a domain U which contains all eigenvalues
 162 $\lambda \in \sigma(T)$ of T . As noted above, implementing an arbitrary holomorphic g is however too costly, and
 163 we instead approximate g via a collection of simpler basis functions $\{\Psi_i\}_{i \in I}$ as $g(\lambda) \approx \sum_{i \in I} \theta_i \Psi_i(\lambda)$.

164 In order to choose the filter bank $\{\Psi_i\}_{i \in I}$, we thus need to answer the question of how to optimally
 165 approximate arbitrary holomorphic functions on a given fixed domain U . The solution to this problem
 166 is given in the guise of **Faber polynomials** [16, 12] which generalize the familiar Chebychev
 167 polynomials utilized in [15] to subsets U of the complex plane [17]. Faber polynomials provide
 168 near near mini-max polynomial approximation³ to any holomorphic function defined on a domain
 169 U satisfying some minimal conditions (c.f. [17] for exact details). What is more, they have already
 170 successfully been employed in numerically approximating matrices of the form $g(T)$ for T not
 171 necessarily symmetric [41].

172 While for a generic domain U Faber polynomials are impossible to compute analytically, this poses
 173 no limitations to us in practice: Short of a costly explicit calculation of the spectrum $\sigma(T)$, the only
 174 information that is generically available, is that eigenvalues may be located anywhere within a circle
 175 of radius $\|T\|$. This circle must thus be contained in any valid domain U . Making the minimal choice
 176 by taking U to be exactly this circle, the n^{th} -Faber polynomial may be analytically calculated [22]:
 177 Up to normalization (absorbed into the learnable parameters) it is given by the monomial λ^n . We
 178 thus take our n^{th} basis element $\Psi_n(\lambda)$ to be given precisely by this monomial: $\Psi_n(\lambda) = \lambda^n$.

179 In a setting where more detailed information on $\sigma(T)$ is available, the domain U may of course be
 180 adapted to reflect this. Corresponding Faber polynomials might then be pre-computed numerically.

181 **3.3.2 Unbounded Spectral Domain: Functions decaying at complex infinity**

182 In the multi-graph setting – e.g. during graph classification – we are confronted with the possibility
 183 that distinct graphs may describe the same underlying object [34, 39, 32]. This might for example
 184 occur if two distinct graphs discretize the same underlying continuous space; e.g. at different
 185 resolution scales. In this setting – instead of precise placements of nodes – what is actually important
 186 is the overall structure and geometry of the respective graphs.

187 Un-normalized Laplacians provide a convenient multi-scale descriptions of such graphs, as they
 188 encode information corresponding to coarse geometry into small (in modulus) eigenvalues, while
 189 finer graph structures correspond to larger eigenvalues [11, 42]. When designing networks whose
 190 outputs are not overly sensitive to fine print articulations of graphs, the spectral response (3) then
 191 provides guidance on determining which holomorphic filters g are able to suppress this superfluous
 192 high-lying spectral information: It is sufficient that $g^{(n)}(\lambda)/n! \approx 0$ for $|\lambda| \gg 1$.

193 It can be shown that no holomorphic function with such large- $|\lambda|$ asymptotics defined on all of \mathbb{C}
 194 exists.⁴ We thus make the minimal necessary change and assume g to be defined on a *punctured*
 195 domain $U = \mathbb{C} \setminus \{y\}$ instead. The choice of $y \in \mathbb{C}$ is treated as a hyperparameter, which may be
 196 adjusted to the task at hand. Any such g may then be expanded as $g(\lambda) = \sum_{j=1}^{\infty} \theta_j (\lambda - y)^{-j}$
 197 for some coefficients $\{\theta_i\}_{i=1}^{\infty}$ [2]. Evaluating the defining integral (2) for the Laplacian L^{in} on the
 198 atoms $\Psi_j(\lambda) = (\lambda - y)^{-j}$ yields $\Psi_j(L^{\text{in}}) = ([L^{\text{in}} - y \cdot Id]^{-1})^j$; as proved in Appendix E. Hence
 199 corresponding filters are polynomials in the **resolvent** $R_y(L^{\text{in}}) := [L^{\text{in}} - y \cdot Id]^{-1}$ of L^{in} .

200 Such resolvents are traditionally used as tools to compare operators with potentially divergent
 201 norms [52]. Recently [33] utilized them in the *undirected* setting to construct networks provably
 202 assigning similar feature-vectors to weighted graphs describing the same underlying object at different
 203 resolution-scales. Our approach extends these networks to the *directed* setting:

³I.e. minimizing the maximal approximation error on the domain of definition U .

⁴This is an immediate consequence of Liouville’s theorem in complex analysis [1].

204 **Effective directed Limit Graphs:** From a diffusion perspective, information in a graph equalizes
 205 much faster along edges with large weights than via weaker edges. In the limit where the edge-weights
 206 within certain sub-graphs tend to infinity, information within these clusters equalizes immediately
 207 and such sub-graphs should thus effectively behave as single nodes. Extending *undirected*-graph
 208 results [33], we here establish rigorously that this is indeed also true in the *directed* setting.

209 Mathematically, we make our arguments precise by considering a graph G with a weight matrix
 210 W admitting a (disjoint) two-scale decomposition as $W = W^{\text{regular}} + c \cdot W^{\text{high}}$ (c.f. Fig. 5). As
 211 the larger weight scale $c \gg 1$ tends to infinity, we then establish that the resolvent $R_y(L^{\text{in}})$ on G
 212 converges to the resolvent $R_y(\underline{L}^{\text{in}})$ of the Laplacian $\underline{L}^{\text{in}}$ on a coarse-grained limit graph \underline{G} . This
 213 limit \underline{G} arises by collapsing the reaches R of the graph $G_{\text{high}} = (G, W^{\text{high}})$ (c.f. Fig. 5 (c)) into
 214 single nodes. For technical reasons, we here assume⁵ equal in- and out-degrees within G_{high} (i.e.
 215 $\sum_i W_{ij}^{\text{high}} = \sum_i W_{ji}^{\text{high}}$). Appendix G contains proofs corresponding to the results below.

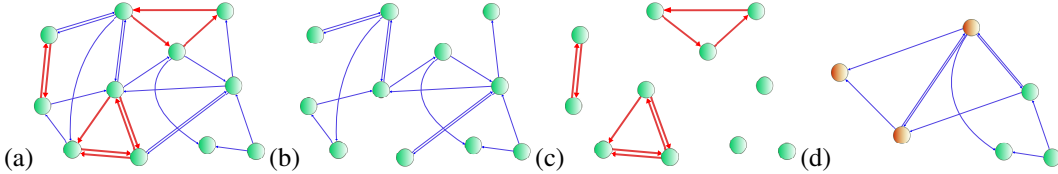


Figure 5: (a) Graph G with W^{regular} (blue) & W^{high} (red); (b) W^{regular} ; (c) W^{high} ; (d) Limit Graph \underline{G}

216 When defining \underline{G} , *directed reaches* now replace the *undirected components* of [33]:

217 **Definition 3.1.** The node set $\underline{\mathcal{G}}$ of \underline{G} is constituted by the set of all reaches in G_{high} . Edges $\underline{\mathcal{E}}$ of
 218 \underline{G} are given by those elements $(R, P) \in \underline{\mathcal{G}} \times \underline{\mathcal{G}}$ with non-zero agglomerated edge weight $\underline{W}_{RP} =$
 219 $\sum_{r \in R} \sum_{p \in P} W_{rp}$. Node weights in \underline{G} are defined similarly by aggregating as $\underline{\mu}_R = \sum_{r \in R} \mu_r$.

220 To map signals between these graphs, translation operators J^\downarrow, J^\uparrow are needed. Let x be a scalar graph
 221 signal and let $\mathbb{1}_R$ be the vector that has 1 as entry for nodes $r \in R$ and zero otherwise. Denote by
 222 u_R the entry of u at node $R \in \underline{\mathcal{G}}$. The projection operator J^\downarrow is then defined component-wise by
 223 evaluation at node $R \in \underline{\mathcal{G}}$ as $(J^\downarrow x)_R = \langle \mathbb{1}_R, x \rangle / \mu_R$. Interpolation is defined as $J^\uparrow u = \sum_{R \in \underline{\mathcal{G}}} u_R \cdot \mathbb{1}_R$.
 224 The maps J^\uparrow, J^\downarrow are then extended from single features $\{x\}$ to feature *matrices* $\{X\}$ via linearity.

225 With these preparations, we can now rigorously establish the suspected effective behaviour of clusters:

226 **Theorem 3.2.** In the above setting, we have $\|R_y(L^{\text{in}}) \cdot X - J^\uparrow R_y(\underline{L}^{\text{in}}) J^\downarrow \cdot X\| \rightarrow 0$ as $c \rightarrow \infty$.

227 For $c \gg 1$, applying the resolvent $R_y(L^{\text{in}})$ on G is thus essentially the same as first projecting
 228 to the coarse-grained graph \underline{G} (where all strongly connected clusters are collapsed), applying the
 229 corresponding resolvent there and then interpolating back up to G . The geometric information within
 230 $R_y(L^{\text{in}})$ is thus essentially reduced to that of the coarse grained geometry within \underline{G} .

231 Large weights within a graph typically correspond to fine-structure articulations of its geometry:
 232 For graph-discretisations of continuous spaces, edge weights e.g. typically correspond to inverse
 233 discretization lengths ($w_{ij} \sim 1/d_{ij}$) and strongly connected clusters describe closely co-located
 234 nodes. In social-networks, edge weights might encode a closeness-measure, and coarse-graining
 235 would correspond to considering interactions between (tightly connected) communities as opposed to
 236 individual users. In either case, fine print articulations are discarded when applying resolvents.

237 **Stability of Filters:** This reduction to a limit description on \underline{G} is respected by our filters $\{g_\theta\}$:

238 **Theorem 3.3.** In the above setting, we have $\|g_\theta(L^{\text{in}}) \cdot X - J^\uparrow g_\theta(\underline{L}^{\text{in}}) J^\downarrow \cdot X\| \rightarrow 0$ as $c \rightarrow \infty$.

239 If the weight-scale c is very large, applying the learned filter $g_\theta(\lambda) = \sum_{i=1}^K \theta_i (\lambda - y)^{-i}$ to a signal
 240 X on G as $X \mapsto g_\theta(L^{\text{in}}) \cdot X$ thus is essentially the same as first discarding fine-structure information
 241 by projecting X to \underline{G} , applying the spectral filter g_θ there and subsequently interpolating back to G .
 242 Information about the precise articulation of a given graph G is thus suppressed in this propagation
 243 scheme; it is purely determined by the graph structure of the coarse-grained description \underline{G} . Theorem
 244 4.2 below establishes that this behaviour persists for entire (directed) spectral convolutional networks.

⁵This is known as Kirchhoff's assumption [4]; reproducing the eponymous law of electrical circuits.

245 **4 Spectral networks on directed graphs: HoloNets**

246 We now collect holomorphic filters into corresponding networks; termed **HoloNets**. In doing so, we
 247 need to account for the possibility that given edge directionalities might limit the information-flow
 248 facilitated by filters $\{g_\theta(T)\}$: In the path-graph setting of Fig. 4 for example, a polynomial filter in the
 249 adjacency matrix would only transport information *along* the graph; features of earlier nodes would
 250 never be augmented with information about later nodes. To circumvent this, we allow for two sets of
 251 filters $\{g_\theta^{\text{fwd}}(T)\}$ and $\{g_\theta^{\text{bwd}}(T^*)\}$ based on the characteristic operator T and its adjoint T^* . Allowing
 252 these forward- and backward-filters to be learned in different bases $\{\Psi_i^{\text{fwd/bwd}}\}_{i \in I^{\text{fwd/bwd}}}$, we may write
 253 $g_\theta^{\text{fwd/bwd}}(\lambda) = \sum_{i \in I^{\text{fwd/bwd}}} \theta_i^{\text{fwd/bwd}} \Psi_i^{\text{fwd/bwd}}(\lambda)$. With bias matrices $B^{\ell+1}$ of size $N \times F_{\ell+1}$ and weight
 254 matrices $W_k^{\text{fwd/bwd}, \ell+1}$ of dimension $F_\ell \times F_{\ell+1}$, our update rule is then efficiently implemented as

$$X^\ell = \rho \left(\alpha \sum_{i \in I^{\text{fwd}}} \Psi_i^{\text{fwd}}(T) \cdot X^{\ell-1} \cdot W_i^{\text{fwd}, \ell} + (1 - \alpha) \sum_{i \in I^{\text{bwd}}} \Psi_i^{\text{bwd}}(T^*) \cdot X^{\ell-1} \cdot W_i^{\text{bwd}, \ell} + B^\ell \right).$$

255 Here ρ is a point-wise non-linearity, and the parameter $\alpha \in [0, 1]$ – learnable or tunable – is
 256 introduced following [47] to allow for a prejudiced weighting of the forward or backward direction.
 257 The generically complex weights & biases may often be restricted to \mathbb{R} without losing expressivity:

258 **Theorem 4.1.** Suppose for filter banks $\{\Psi_i^{\text{fwd/fwd}}\}_{I^{\text{fwd/fwd}}}$ that the matrices $\Psi_i^{\text{fwd}}(T)$, $\Psi_i^{\text{bwd}}(T^*)$ contain
 259 only real entries. Then any HoloNet with layer-widths $\{F_\ell\}$ with complex weights & biases and a non
 260 linearity that acts on complex numbers componentwise as $\rho(a + ib) = \tilde{\rho}(a) + i\tilde{\rho}(b)$ can be exactly
 261 represented by a HoloNet of widths $\{2 \cdot F_\ell\}$ utilizing $\tilde{\rho}$ and employing only real weights & biases.

262 This result (proved in Appendix H) establishes that for the same number of *real* parameters, real
 263 HoloNets theoretically have greater expressive power than complex ones. In our experiments in
 264 Section 5 below, we empirically find that complex weights do provide advantages on some graphs.
 265 Thus we propose to treat the choice of complex vs. real parameters as a binary hyperparameter.

266 **FaberNet:** The first specific instantiation of HoloNets we consider, employs the Faber Polynomials
 267 of Section 3.3.1 for both the forward and backward filter banks. Since [47] established that considering
 268 edge directionality is especially beneficial on heterophilic graphs, this is also our envisioned target
 269 for the corresponding networks. We thus use as characteristic operator a matrix that avoids direct
 270 comparison of feature vectors of a node with those of immediate neighbours: We choose $T =$
 271 $(D^{\text{in}})^{-\frac{1}{4}} \cdot W \cdot (D^{\text{out}})^{-\frac{1}{4}}$ since it has a zero-diagonal and its normalization performed well empirically.
 272 For the same reason of heterophily, we also consider the choice of whether to include the Faber
 273 polynomial $\Psi_0(\lambda) = 1$ in our basis as a hyperparameter. As non-linearity, we choose either
 274 $\rho(a + ib) = \text{ReLU}(a) + i\text{ReLU}(b)$ or $\rho(a + ib) = |a| + i|b|$. Appendix I contains additional details.

275 **Dir-ResolvNet:** In order to build networks that are insensitive to the fine-print articulation of
 276 *directed* graphs, we take as filter bank the functions $\{\Psi_i(\lambda) = (\lambda - y)^{-i}\}_{i>0}$ evaluated on the
 277 Laplacian L^{in} for both the forward and backward direction. To account for individual node-weights
 278 when building up graph-level features, we use an aggregation Ω that aggregates $N \times F$ -dimensional
 279 node-feature matrices as $\Omega(X)_j = \sum_{i=1}^N |X_{ij}| \cdot \mu_i$ to a graph-feature $\Omega(X) \in \mathbb{R}^F$. Graph-level
 280 stability under varying resolution scales is then captured by our next result:

281 **Theorem 4.2.** Let Φ and $\underline{\Phi}$ be the feature maps associated to Dir-ResolvNets with the same weights
 282 and biases deployed on graphs G and \underline{G} as defined in Section 3.3.2. With Ω the aggregation method
 283 specified above and $W = W^{\text{regular}} + c \cdot W^{\text{high}}$ as in Theorem 3.3, we have for $c \rightarrow \infty$:

$$\|\Omega(\Phi(X)) - \Omega(\underline{\Phi}(J^\downarrow X))\| \longrightarrow 0$$

284 Appendix G contains proofs of this and additional stability results. From Theorem 4.2 we conclude
 285 that graph-level features generated by a Dir-ResolvNet are indeed insensitive to fine print articulations
 286 of weighted digraphs: As discussed in Section 3.3.2, geometric information corresponding to such
 287 fine details is typically encoded into strongly connected sub-graphs; with the connection strength c
 288 corresponding to the level of detail. However, information about the structure of these sub-graphs is
 289 precisely what is discarded when moving to \underline{G} via J^\downarrow . Thus the greater the level of detail within G ,
 290 the more similar are generated feature-vectors to those of a (relatively) coarse-grained description \underline{G} .
 291

292 5 Experiments

293 We present experiments on real-world data to evaluate the capabilities of our HoloNets numerically.

294 5.1 FaberNet: Node Classification

295 We first evaluate on the task of node-classification in the presence of heterophily. We consider multiple
 296 heterophilic graph-datasets on which we compare the performance of our **FaberNet** instantiation of
 297 the HoloNet framework against a representative array of baselines: As *simple baselines* we consider
 298 MLP and GCN [31]. H₂GCN [60], GPR-GNN [10], LINKX [37], FSGNN [40], ACM-GCN [38],
 299 GloGNN [36] and Gradient Gating [51] constitute *heterophilic state-of-the-art models*. Finally *state-*
 300 *of-the-art models for directed graphs* are given by DiGCN [53], MagNet [59] and Dir-GNN [47].
 301 Appendix I contains dataset statistics as well as additional details on baselines, experimental setup
 302 and hyperparameters.

Table 1: Results on real-world directed heterophilic datasets. OOM indicates out of memory.

	Squirrel	Chameleon	Arxiv-year	Snap-patents	Roman-Empire
Homophily	0.223	0.235	0.221	0.218	0.05
MLP	28.77 ± 1.56	46.21 ± 2.99	36.70 ± 0.21	31.34 ± 0.05	64.94 ± 0.62
GCN	53.43 ± 2.01	64.82 ± 2.24	46.02 ± 0.26	51.02 ± 0.06	73.69 ± 0.74
H ₂ GCN	37.90 ± 2.02	59.39 ± 1.98	49.09 ± 0.10	OOM	60.11 ± 0.52
GPR-GNN	54.35 ± 0.87	62.85 ± 2.90	45.07 ± 0.21	40.19 ± 0.03	64.85 ± 0.27
LINKX	61.81 ± 1.80	68.42 ± 1.38	56.00 ± 0.17	61.95 ± 0.12	37.55 ± 0.36
FSGNN	74.10 ± 1.89	78.27 ± 1.28	50.47 ± 0.21	65.07 ± 0.03	79.92 ± 0.56
ACM-GCN	67.40 ± 2.21	74.76 ± 2.20	47.37 ± 0.59	55.14 ± 0.16	69.66 ± 0.62
GloGNN	57.88 ± 1.76	71.21 ± 1.84	54.79 ± 0.25	62.09 ± 0.27	59.63 ± 0.69
Grad. Gating	64.26 ± 2.38	71.40 ± 2.38	63.30 ± 1.84	69.50 ± 0.39	82.16 ± 0.78
DiGCN	37.74 ± 1.54	52.24 ± 3.65	OOM	OOM	52.71 ± 0.32
MagNet	39.01 ± 1.93	58.22 ± 2.87	60.29 ± 0.27	OOM	88.07 ± 0.27
DirGNN	75.13 ± 1.95	79.74 ± 1.40	63.97 ± 0.30	73.95 ± 0.05	91.3 ± 0.46
FaberNet	76.71 ± 1.92	80.33 ± 1.19	64.62 ± 1.01	75.10 ± 0.03	92.24 ± 0.43

303 As can be inferred from Table 1, **FaberNet sets new state of the art results** on all five heterophilic
 304 graph datasets above; out-performing intricate *undirected* methods specifically designed for the
 305 setting of heterophily. What is more, it also significantly outperforms *directed spatial* methods such
 306 as Dir-GNN, whose results can be considered as reporting a best-of-three performance over multiple
 307 directed spatial methods (c.f. Appendix I or [47] for details). FaberNet also significantly out-performs
 308 MagNet. This method is a spectral model, which relies on the graph Fourier transform associated
 309 to a certain operator that is able to remain self-adjoint in the directed setting. We thus might take
 310 this gap in performance as further evidence of the utility of transcending the classical graph Fourier
 311 transform: Utilizing the holomorphic functional calculus – as opposed to the traditional graph Fourier
 312 transform – allows to base filters on (non-self-adjoint) operators more adapted to the respective task at
 313 hand. On Squirrel and Chameleon, our method performed best when using complex parameters (c.f.
 314 Table 7 in Appendix I). With MagNet being the only other method utilizing complex parameters, its
 315 performance gap to Dir-ResolvNet also implies that it is indeed the interplay of complex weights with
 316 judiciously chosen filter banks and characteristic operators that provides state-of-the-art performance;
 317 not the use of complex parameters alone.

318 5.2 Dir-ResolvNet: DiGraph Regression and Scale-Insensitivity

319 We test the properties of our **Dir-ResolvNet** HoloNet via graph regression experiments. Weighted-
 320 directed datasets containing both node-features and graph-level targets are currently still scarce. Hence
 321 we follow [33] and evaluate on the task of molecular property prediction. While neither our Dir-
 322 ResolvNet nor baselines of Table 1 are designed for this task, such molecular data still allows for fair
 323 comparisons of expressive power and stability properties of (non-specialized) graph learning methods
 324 [28]. We utilize the QM7 dataset [50], containing graphs of 7165 organic molecules; each containing
 325 hydrogen and up to seven types of heavy atoms. Prediction target is the molecular atomization energy.
 326 While each molecule is originally represented by its Coulomb matrix $W_{ij} = Z_i \cdot Z_j / |\vec{x}_i - \vec{x}_j|$, we

327 modify these edge-weights: Between each heavy atom and all atoms outside its respective immediate
 328 hydrogen cloud we set $W_{ij} = Z_i^{\text{outside}} \cdot (Z_j^{\text{heavy}} - 1) / |\vec{x}_i - \vec{x}_j|$. While the sole reason for this change
 329 is to make the underlying graphs directed (enabling comparisons of *directed* methods), we
 330 might heuristically interpret it as arising from a (partial) shielding of heavy atoms by surrounding
 331 electrons [26].

332 **Digraph-Regression:** Treating W as a directed weight-matrix, we evaluate Dir-ResolvNet against
 333 all other *directed* methods of Table 1. Atomic charges are used as node weights ($\mu_i = Z_i$) where
 applicable and one-hot encodings of atomic charges Z_i provide node-wise
 input features. As evident from Table 2, our method produces significantly
 334 lower mean-absolute-errors (MAEs) than corresponding directed baselines
 (c.f. Table 1): **Competitors are out-performed by a factor of two and
 more.** We attribute this to Dir-ResolvNets ability to discard superfluous
 information; thus better representing overall molecular geometries.

DirGNN	59.01±2.54
MagNet	45.31±4.24
DiGCN	39.95±6.23
Dir-ResolvNet	17.12±0.63

335 **Scale-Insensitivity:** To numerically investigate the stability properties of Dir-ResolvNet that were

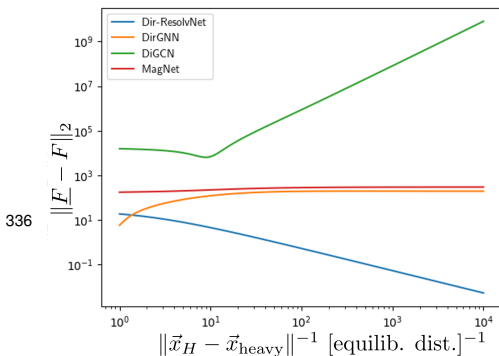


Figure 6: Feature-difference for collapsed (F_c) and deformed (F) graphs.

337 In this setting, we compare feature vectors of collapsed graphs with feature vectors of molecules
 338 where hydrogen atoms have been deflected but have not yet arrived at the positions of nearest heavy
 339 atoms. Feature vectors are generated with the previously trained networks of Table 2. As evident
 340 from Fig. 6, **Dir-ResolvNet’s feature-vectors converge** as the scale $c \sim \|\vec{x}_H - \vec{x}_{\text{heavy}}\|^{-1}$ increases;
 341 thus numerically verifying the scale-invariance Theorem 4.2. **Feature vectors of baselines do not
 342 converge:** These models are sensitive to changes in resolutions when generating graph-level features.

343 This difference in sensitivity is also apparent in our final experiment, where collapsed molecular
 344 graphs $\{\underline{G}\}$ are treated as a model for data obtained from a resolution-limited observation process un-
 345 able to resolve individual H-atoms. Given models trained on directed higher resolution digraphs $\{G\}$,
 atomization energies are then to be predicted solely using coarse
 grained molecular digraphs. While Dir-ResolvNet’s prediction accu-
 346 racy remains high, performance of baselines decreases significantly if
 the resolution scale is reduced during inference: While Dir-ResolvNet
 out-performed baselines by a factor of two and higher before, this **lead
 increases to a factor of up to 240 if resolutions vary** (c.f. Table 3).

DirGNN	195.64±2.20
MagNet	663.63±190.358
DiGCN	6672.71±2243.61
Dir-ResolvNet	27.34±7.55

347 6 Conclusion

348 We introduced the HoloNet framework, which allows to extend spectral networks to directed graphs.
 349 Key building blocks of these novel networks are newly introduced holomorphic filters, no longer
 350 reliant on the graph Fourier transform. We provided a corresponding frequency perspective, inves-
 351 tigated optimal filter-banks and discussed the interplay of filters with characteristic operators in
 352 shaping inductive biases. Experiments on real world data considered two particular HoloNet instan-
 353 tiations: FaberNet provided new state-of-the-art results for node classification under heterophily while
 354 Dir-ResolvNet generated feature vectors stable to resolution-scale-varying topological perturbations.

References

- [1] L. V. Ahlfors. *Complex Analysis*. McGraw-Hill Book Company, 2 edition, 1966.
- [2] Joseph Bak and Donald J. Newman. *Complex analysis*. Springer, 2017.
- [3] Muhammet Balcilar, Guillaume Renton, Pierre Héroux, Benoit Gaüzère, Sébastien Adam, and Paul Honeine. Analyzing the expressive power of graph neural networks in a spectral perspective. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [4] Marwa Balti. Non self-adjoint laplacians on a directed graph, 2018.
- [5] Dominique Beaini, Saro Passaro, Vincent Létourneau, William L. Hamilton, Gabriele Corso, and Pietro Lió. Directional graph networks. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 748–758. PMLR, 2021.
- [6] Filippo Maria Bianchi, Daniele Grattarola, Lorenzo Francesco Livi, and Cesare Alippi. Graph neural networks with convolutional arma filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44:3496–3507, 2019.
- [7] L. C. Blum and J.-L. Reymond. 970 million druglike small molecules for virtual screening in the chemical universe database GDB-13. *J. Am. Chem. Soc.*, 131:8732, 2009.
- [8] Deyu Bo, Chuan Shi, Lele Wang, and Renjie Liao. Specformer: Spectral graph neural networks meet transformers. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [9] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann Lecun. Spectral networks and locally connected networks on graphs. In *International Conference on Learning Representations (ICLR2014), CBLS, April 2014*, 2014.
- [10] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank graph neural network. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [11] Fan R. K. Chung. *Spectral graph theory*, volume 92 of *CBMS Regional Conference Series in Mathematics*. Conference Board of the Mathematical Sciences, Washington, DC; by the American Mathematical Society, Providence, RI, 1997.
- [12] John P. Coleman and Russell A. Smith. The faber polynomials for circular sectors. *Mathematics of Computation*, 49(179):231–241, 1987.
- [13] Fabrizio Colombo, Irene Sabadini, and Daniele C. Struppa. *Noncommutative Functional Calculus: Theory and Applications of Slice Hyperholomorphic Functions*, pages 201–210. Springer Basel, Basel, 2011.
- [14] Michael Cowling, Ian Doust, Alan Micintosh, and Atsushi Yagi. Banach space operators with a bounded hoo functional calculus. *Journal of the Australian Mathematical Society*, 60(1):51–89, 1996.
- [15] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29, 2016.
- [16] S. W. Ellacott. Computation of faber series with application to numerical polynomial approximation in the complex plane. *Mathematics of Computation*, 40(162):575–587, 1983.
- [17] G. H. Elliott. The construction of chebyshev approximations in the complex plane. PhD Thesis. Imperial College London, 1978.
- [18] Simon Geisler, Yujia Li, Daniel J. Mankowitz, Ali Taylan Cemgil, Stephan Günnemann, and Cosmin Paduraru. Transformers meet directed graphs. *CoRR*, abs/2302.00049, 2023.

- 402 [19] Herbert A. Gindler. An operational calculus for meromorphic functions. *Nagoya Mathematical*
403 *Journal*, 26:31–38, 1966.
- 404 [20] Markus Haase. *The Functional Calculus for Sectorial Operators*, pages 19–60. Birkhäuser
405 Basel, Basel, 2006.
- 406 [21] William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on
407 large graphs. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob
408 Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information*
409 *Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017,*
410 *December 4-9, 2017, Long Beach, CA, USA*, pages 1024–1034, 2017.
- 411 [22] M. He. The faber polynomials for circular lunes. *Computers & Mathematics with Applications*,
412 30(3):307–315, 1995.
- 413 [23] Mingguo He, Zhewei Wei, Zengfeng Huang, and Hongteng Xu. Bernnet: Learning arbitrary
414 graph spectral filters via bernstein approximation. In Marc’Aurelio Ranzato, Alina Beygelzimer,
415 Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural*
416 *Information Processing Systems 34: Annual Conference on Neural Information Processing*
417 *Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 14239–14251, 2021.
- 418 [24] Mingguo He, Zhewei Wei, and Ji-Rong Wen. Convolutional neural networks on graphs with
419 chebyshev approximation, revisited. In *NeurIPS*, 2022.
- 420 [25] Yixuan He, Michael Perlmutter, Gesine Reinert, and Mihai Cucuringu. MSGNN: A spectral
421 graph neural network based on a novel magnetic signed laplacian. In Bastian Rieck and Razvan
422 Pascanu, editors, *Learning on Graphs Conference, LoG 2022, 9-12 December 2022, Virtual*
423 *Event*, volume 198 of *Proceedings of Machine Learning Research*, page 40. PMLR, 2022.
- 424 [26] Hendrik Heinz and Ulrich W. Suter. Atomic charges for classical simulations of polar systems.
425 *The Journal of Physical Chemistry B*, 108(47):18341–18352, 2004.
- 426 [27] Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- 427 [28] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
428 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.
429 In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-
430 Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference*
431 *on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual,*
432 2020.
- 433 [29] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
434 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.
435 *arXiv preprint arXiv:2005.00687*, 2020.
- 436 [30] Tosio Kato. *Perturbation theory for linear operators; 2nd ed.* Grundlehren der mathematischen
437 Wissenschaften : a series of comprehensive studies in mathematics. Springer, Berlin, 1976.
- 438 [31] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
439 networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon,*
440 *France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- 441 [32] Christian Koke. Limitless stability for graph convolutional networks. In *11th International*
442 *Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenRe-
443 view.net, 2023.
- 444 [33] Christian Koke, Abhishek Saroha, Yuesong Shen, Marvin Eisenberger, and Daniel Cremers.
445 Resolvnet: a graph convolutional network with multi-scale consistency, 2023.
- 446 [34] Ron Levie, Michael M. Bronstein, and Gitta Kutyniok. Transferability of spectral graph
447 convolutional neural networks. *CoRR*, abs/1907.12972, 2019.
- 448 [35] Ron Levie, Federico Monti, Xavier Bresson, and Michael M. Bronstein. Cayleynets: Graph
449 convolutional neural networks with complex rational spectral filters. *IEEE Trans. Signal*
450 *Process.*, 67(1):97–109, 2019.

- 451 [36] Xiang Li, Renyu Zhu, Yao Cheng, Caihua Shan, Siqiang Luo, Dongsheng Li, and Weining
452 Qian. Finding global homophily in graph neural networks when meeting heterophily. In
453 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan
454 Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022,*
455 *Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages
456 13242–13256. PMLR, 2022.
- 457 [37] Derek Lim, Felix Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Bhalerao,
458 and Ser-Nam Lim. Large scale learning on non-homophilous graphs: New benchmarks and
459 strong simple methods. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy
460 Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing*
461 *Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS*
462 *2021, December 6-14, 2021, virtual*, pages 20887–20902, 2021.
- 463 [38] Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen
464 Chang, and Doina Precup. Revisiting heterophily for graph neural networks. In *NeurIPS*, 2022.
- 465 [39] Sohir Maskey, Ron Levie, and Gitta Kutyniok. Transferability of graph neural networks: an
466 extended graphon approach. *CoRR*, abs/2109.10096, 2021.
- 467 [40] Sunil Kumar Maurya, Xin Liu, and Tsuyoshi Murata. Improving graph neural networks with
468 simple architecture design. *CoRR*, abs/2105.07634, 2021.
- 469 [41] I. Moret and P. Novati. The computation of functions of matrices by truncated faber series.
470 *Numerical Functional Analysis and Optimization*, 22(5-6):697–719, 2001.
- 471 [42] Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm.
472 In T. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information*
473 *Processing Systems*, volume 14. MIT Press, 2001.
- 474 [43] Hoang NT and Takatori Maehara. Revisiting graph neural networks: All we have is low-pass
475 filters. *CoRR*, abs/1905.09550, 2019.
- 476 [44] Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn:
477 Geometric graph convolutional networks, 2020.
- 478 [45] Oleg Platonov, Denis Kuznedelev, Michael Diskin, Artem Babenko, and Liudmila
479 Prokhorenkova. A critical look at the evaluation of gnns under heterophily: Are we really
480 making progress? In *The Eleventh International Conference on Learning Representations,*
481 *ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- 482 [46] Olaf Post. *Spectral Analysis on Graph-like Spaces / by Olaf Post*. Lecture Notes in Mathematics,
483 2039. Springer Berlin Heidelberg, Berlin, Heidelberg, 1st ed. 2012. edition, 2012.
- 484 [47] Emanuele Rossi, Bertrand Charpentier, Francesco Di Giovanni, Fabrizio Frasca, Stephan
485 Günnemann, and Michael Bronstein. Edge directionality improves learning on heterophilic
486 graphs, 2023.
- 487 [48] Luana Ruiz, Luiz F. O. Chamon, and Alejandro Ribeiro. Transferability properties of graph
488 neural networks. *CoRR*, abs/2112.04629, 2021.
- 489 [49] Luana Ruiz, Fernando Gama, and Alejandro Ribeiro. Graph neural networks: Architectures,
490 stability and transferability, 2021.
- 491 [50] M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld. Fast and accurate modeling
492 of molecular atomization energies with machine learning. *Physical Review Letters*, 108:058301,
493 2012.
- 494 [51] T. Konstantin Rusch, Benjamin Paul Chamberlain, Michael W. Mahoney, Michael M. Bronstein,
495 and Siddhartha Mishra. Gradient gating for deep multi-rate learning on graphs. In *The Eleventh*
496 *International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5,*
497 *2023*. OpenReview.net, 2023.

- 498 [52] Gerald Teschl. *Mathematical Methods in Quantum Mechanics*. American Mathematical Society,
499 2014.
- 500 [53] Zekun Tong, Yuxuan Liang, Changsheng Sun, Xinke Li, David Rosenblum, and Andrew Lim.
501 Digraph inception convolutional networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F.
502 Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,
503 pages 17907–17918. Curran Associates, Inc., 2020.
- 504 [54] J. J. P. Veerman and Robert Lyons. A primer on laplacian dynamics in directed graphs, 2020.
- 505 [55] Jjp Veerman and Ewan Kummel. Diffusion and consensus on weakly connected directed graphs.
506 *Linear Algebra and its Applications*, 578, 05 2019.
- 507 [56] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
508 Bengio. Graph attention networks. In *6th International Conference on Learning Representations,*
509 *ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings.*
510 OpenReview.net, 2018.
- 511 [57] Xiyuan Wang and Muhan Zhang. How powerful are spectral graph neural networks. In
512 *International Conference on Machine Learning*, 2022.
- 513 [58] Xiyuan Wang and Muhan Zhang. How powerful are spectral graph neural networks. In
514 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan
515 Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022,*
516 *Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages
517 23341–23362. PMLR, 2022.
- 518 [59] Xitong Zhang, Yixuan He, Nathan Brugnone, Michael Perlmutter, and Matthew J. Hirn. Magnet:
519 A neural network for directed graphs. In Marc’ Aurelio Ranzato, Alina Beygelzimer, Yann N.
520 Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information*
521 *Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021,*
522 *NeurIPS 2021, December 6-14, 2021, virtual*, pages 27003–27015, 2021.
- 523 [60] Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra.
524 Beyond homophily in graph neural networks: Current limitations and effective designs. In Hugo
525 Larochelle, Marc’ Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin,
526 editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural*
527 *Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.

528 **A Notation**

529 We provide a summary of employed notational conventions:

Table 4: Classification Accuracies on Social Network Datasets

Symbol	Meaning
G	a graph
\mathcal{G}	Nodes of the graph G
N	number of nodes $ \mathcal{G} $ in G
\underline{G}	Coarse grained version of graph G
μ_i	weight of node i
M	weight matrix
$\langle \cdot, \cdot \rangle$	inner product
W	(weighted) adjacency matrix
$D^{\text{in/out}}$	in/out-degree matrix
L^{in}	in-degree graph Laplacian
T	generic characteristic operator
T^*	hermitian adjoint of T
T^\top	transpose of T (used if and only if T has only real entries)
U	change-of-basis matrix to a (complete) basis consisting of eigenvectors (used in the undirected setting only)
$\sigma(T)$	spectrum (i.e. collection of eigenvalues) of T
λ	an eigenvalue
g	a holomorphic function
$g(T)$	function g applied to operator T
Ψ_i	an element of a filter-bank
z, y	complex numbers
U	subset of \mathbb{C}
$R_z(L^{\text{in}})$	the resolvent $(L^{\text{in}} - z \cdot Id)^{-1}$
c	a weight-/resolution- scale
J^\downarrow, J^\uparrow	projection and interpolation operator
Φ	map associated to a graph convolution network
Ω	graph-level aggregation mechanism
Z_i	atomic charge of atom corresponding to node i
\vec{x}_i	Cartesian position of atom corresponding to node i
$\frac{Z_i Z_j}{ \vec{x}_i - \vec{x}_j }$	Coulomb interaction between atoms i and j
$ \vec{x}_i - \vec{x}_j $	Euclidean distance between x_i and x_j

530 **B Operators beyond the self-adjoint setting**

531 Here we briefly recapitulate facts from linear algebra regarding self-adjoint and non-self-adjoint
 532 matrices, their spectral theory and canonical decompositions.

533 **Self-Adjoint Matrices:** Let us begin with the familiar self adjoint matrices. Given a vector space
 534 V with inner product $\langle \cdot, \cdot \rangle_V$ consider a linear map $T : V \rightarrow V$. The adjoint T^* of the map T is
 535 defined by the demanding that for all $x, y \in V$ we have

$$\langle x, Ty \rangle_V = \langle T^*x, y \rangle_V.$$

536 If we have $V = \mathbb{C}^d$ and the inner product is simply the standard scalar product $\langle x, y \rangle_{\mathbb{C}^d} = \bar{x}^\top y$, we
 537 have

$$T^* = \overline{T}^\top.$$

538 Here \bar{x} denotes the complex conjugate of x and A^\top denotes the transpose of A .

539 An eigenvector- eigenvalue pair, is a pair of a scalar $\lambda \in \mathbb{C}$ and vector $v \in V$ so that

$$(T - \lambda \cdot Id)v = 0.$$

540 It is a standard result in linear algebra (see e.g. [30]) that the spectrum $\sigma(T)$ of all eigenvalues λ of T
 541 contains only real numbers (i.e. $\lambda \in \mathbb{R}$), if T is self-adjoint. What is more, there exists a complete
 542 basis $\{v_i\}_{i=1}^d$ of such eigenvectors that span all of V . These eigenvectors may be chosen to satisfy

$$\langle v_i, v_j \rangle = \delta_{ij}, \quad (4)$$

543 with the Kronecker delta δ_{ij} . As a consequence of this fact, there exists a family of (so called) spectral
 544 projections $\{P_\lambda\}_\lambda$ that have the following properties (c.f. e.g. [52] for a proof)

- 545 • Each P_λ is a linear map on V ($P_\lambda : V \rightarrow V$).
- 546 • For all eigenvectors v to the eigenvalue λ (i.e. v such that $(T - \lambda Id)v = 0$), we have

$$P_\lambda v = v.$$

- 547 • If w is an eigenvector to a different eigenvalue $\mu \in \sigma(T)$ (i.e. w such that $(T - \mu \cdot Id) \cdot w = 0$
 548 and $\mu \neq \lambda$), we have

$$P_\lambda w = 0.$$

- 549 • Each P_λ is a projection (i.e. satisfies $P_\lambda \cdot P_\lambda = P_\lambda$).
- 550 • Each P_λ is self-adjoint (i.e. $P_\lambda = P_\lambda^*$).
- 551 • Each P_λ commutes with T (i.e. $P_\lambda \cdot T = T \cdot P_\lambda$).
- 552 • The family $\{P_\lambda\}_{\lambda \in \sigma(T)}$ form a resolution of the identity:

$$\sum_{\lambda \in \sigma(T)} P_\lambda = Id.$$

553 These properties together then allow us to "diagonalise" the operator T and write it as a sum over its
 554 eigenvalues:

$$T = \sum_{\lambda \in \sigma(T)} \lambda \cdot P_\lambda.$$

555 Applying a generic function may then be defined as [52]

$$g(T) := \sum_{\lambda \in \sigma(T)} g(\lambda) \cdot P_\lambda.$$

556 **Non self-adjoint operators:** If the operator T is no longer self adjoint, eigenvalues no longer need
 557 to be in real, but are generically complex (i.e. $\lambda \in \mathbb{C}$). Furthermore, while there still exist eigenvectors
 558 these no longer need to form a basis of the space V . What instead becomes important in this setting,
 559 are so called *generalized eigenvectors*. A generalized eigenvector w associated to the eigenvalue λ
 560 is a vector for which we have

$$(T - \lambda \cdot Id)^n \cdot w = 0$$

561 for some power $n \in \mathbb{N}$. Clearly each actual *eigenvector* is also a generalized eigenvector (simply for
 562 $n = 1$). What can be shown [30] is that there is a basis of V consisting purely of *generalized* eigen-
 563 vectors (each associated to some eigenvalue λ). These now however need no longer be orthogonal
 564 (i.e. they need not satisfy (4)).

565 There then exists a family of spectral projections $\{P_\lambda\}_\lambda$ that have the following modified set of
 566 properties (c.f. e.g. [30] for a proof)

- 567 • Each P_λ is a linear map on V ($P_\lambda : V \rightarrow V$).
- 568 • For all generalized eigenvectors w to the eigenvalue λ (i.e. w such that $(T - \lambda Id)^n \cdot w = 0$
 569 for some $n \in \mathbb{N}$), we have

$$P_\lambda w = w.$$

- 570 • If w is a generalized eigenvector to a different eigenvalue $\mu \in \sigma(T)$ (i.e. w such that
 571 $(T - \mu \cdot Id)^m \cdot w = 0$ for some $m \in \mathbb{N}$ and $\mu \neq \lambda$), we have

$$P_\lambda w = 0.$$

- 572 • Each P_λ is a projection (i.e. satisfies $P_\lambda \cdot P_\lambda = P_\lambda$).

- 573 • Each P_λ commutes with T (i.e. $P_\lambda \cdot T = T \cdot P_\lambda$).
 574 • The family $\{P_\lambda\}_{\lambda \in \sigma(T)}$ form a resolution of the identity:

$$\sum_{\lambda \in \sigma(T)} P_\lambda = Id.$$

575 Using these facts, we may thus decompose each operator T into a sum as

$$T = \sum_{\lambda \in \sigma(T)} \lambda \cdot P_\lambda + \sum_{\lambda \in \sigma(T)} (T - \lambda \cdot Id) \cdot P_\lambda. \quad (5)$$

576 This sum decomposition is referred to as Jordan-Chevalley decomposition. Importantly, for each
 577 $\lambda \in \sigma(T)$ there is a corresponding natural number $m_\lambda \in \mathbb{N}$ referred to the *algebraic multiplicity* of the
 578 eigenvalue λ . This number m_λ counts, how many generalized eigenvectors $\{w_i^\lambda\}_{i=1}^{m_\lambda}$ are associated to
 579 the generalized eigenvalue λ . These generalized eigenvectors can be chosen such that

$$(T - \lambda \cdot Id) \cdot w_i^\lambda = (T - \lambda \cdot Id) \cdot P_\lambda \cdot w_i^\lambda = w_{i+1}^\lambda,$$

580 if $i \leq m_\lambda - 1$. For the case $i \geq m_\lambda$ we have

$$(T - \lambda \cdot Id) \cdot w_{m_\lambda}^\lambda = 0.$$

581 In total, this implies the nilpotency-relation

$$(T - \lambda \cdot Id)^{m_\lambda} \cdot P_\lambda = [(T - \lambda \cdot Id) \cdot P_\lambda]^{m_\lambda} = 0.$$

582 C Spectrum of adjacency of three-node directed path graph:

583 The adjacency matrix associated to the graph depicted in Figure 4 is given by

$$W = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}.$$

584 Denote by e_i the i^{th} canonical basis vector. Then

$$W \cdot e_1 = e_2, \quad W \cdot e_2 = e_3, \quad W \cdot e_3 = 0.$$

585 Clearly e_3 is an eigenvector to the eigenvalue $\lambda = 0$, so that $0 \in \sigma(W)$. Suppose there exists an
 586 eigenvector v associated to a different eigenvalue $\mu \neq 0$:

$$W \cdot v = \mu \cdot v.$$

587 Then we have

$$W^3 \cdot v = \mu^3 \cdot v$$

588 Since $\{e_1, e_2, e_3\}$ clearly is a basis, there are numbers a, b, c so that

$$v = ae_1 + be_2 + ce_3.$$

589 But then

$$\mu^3 \cdot v = W^3 \cdot v = aW^3 \cdot e_1 + bW^3 \cdot e_2 + cW^3 \cdot e_3 = a \cdot 0 + b \cdot 0 + c \cdot 0.$$

590 Thus we have $\mu = 0$ and hence zero is indeed the only eigenvalue.

591 D Complex Differentiability

592 Complete introductions into this subject may be found in [1] or [2].

593 For a complex valued function $f : \mathbb{C} \rightarrow \mathbb{C}$ of a single complex variable, the derivative of f at a point
 594 $z_0 \in \mathbb{C}$ in its domain of definition is defined as the limit

$$f'(z_0) := \lim_{z \rightarrow z_0} \frac{f(z) - f(z_0)}{z - z_0}.$$

595 For this limit to exist, it needs to be independent of the 'direction' in which z approaches z_0 , which is
 596 a stronger requirement than being real-differentiable.

597 A function is called holomorphic on an open set U if it is complex differentiable at every point in U .
 598 The value $g(\lambda)$ of any such function g at λ may be reproduced as

$$g(\lambda) = -\frac{1}{2\pi i} \oint_S \frac{g(z)}{\lambda - z} dz \quad (6)$$

599 for any circle $S \subseteq \mathbb{C}$ encircling λ so that S is completely contained within U and may be contracted
 600 to a point without leaving U . This equation is referred to as Cauchy's integral formula [2].

601 In fact, the integration contour need not be a circle S , but may be the boundary of any so called
 602 Cauchy domain containing λ :

603 **Definition D.1.** A subset D of the complex plane \mathbb{C} is called a Cauchy domain if D is open, has a
 604 finite number of components (the closure of two of which are disjoint) and the boundary of ∂D of D
 605 is composed of a finite number of closed rectifiable Jordan curves, no two of which intersect.

606 Integrating around any such boundary then reproduces the value of g at λ .

607 E Additional Details on the Holomorphic Functional Calculus

608 **Fundamental Definition:** In order to define the matrix $g(T)$, the formal replacement $\lambda \mapsto T$ is
 609 made on both sides of the Cauchy formula (6), with the path Γ now not only encircling a single value
 610 λ but all eigenvalues $\lambda \in \sigma(T)$ (c.f. also Fig. 7):

$$g(T) := -\frac{1}{2\pi i} \oint_{\Gamma} g(z) \cdot (T - z \cdot Id)^{-1} dz \quad (7)$$

611 The integral is well defined since all eigenvalues of T are assumed to
 lie *inside* the path Γ . For any choice of integration variable z *on* this
 path Γ , the matrix $(T - z \cdot Id)$ is thus indeed invertible, since z is
 never an eigenvalue.

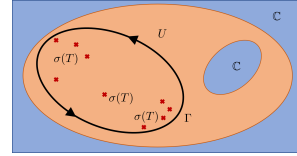


Figure 7: Operator Integral (7)

612 This integral is also known as Dunford-Taylor integral, and the holomorphic functional calculus is
 613 also sometimes called Riesz–Dunford functional calculus [30, 46]. As with the scalar valued integral
 614 (6), it can be shown that the precise path Γ is not important, as long as it encircles all eigenvalues
 615 counter-clockwise and is contractable to a single point within the domain of definition U for the
 616 function g .

617 **Spectral characterization of the Holomorphic Functional Calculus:** As stated in Section 3.2,
 618 the operator $g(T)$ may also be characterized spectrally. Writing $T = \sum_{\lambda \in \sigma(T)} \lambda \cdot P_{\lambda} + \sum_{\lambda \in \sigma(T)} (T -$
 619 $\lambda \cdot Id) \cdot P_{\lambda}$ as in (5), it can then be shown [30], that the spectral action of a given function g is given
 620 as

$$g(T) = \sum_{\lambda \in \sigma(T)} g(\lambda) P_{\lambda} + \sum_{\lambda \in \sigma(T)} \left[\sum_{n=1}^{m_{\lambda}-1} \frac{g^{(n)}(\lambda)}{n!} (T - \lambda \cdot Id)^n \right] P_{\lambda}. \quad (8)$$

621 A proof of this can be found in Chapter 1 of [30].

622 One of the strong properties of complex differentiability previously alluded to is, that as soon as
 623 a function is complex-differentiable once, it is already complex differentiable infinitely often [2].
 624 Hence the n^{th} -derivative in (8) does indeed exist.

625 With the preparations of Appendix B, we can interpret the sum

$$S = \sum_{\lambda \in \sigma(T)} \left[\sum_{n=1}^{m_{\lambda}-1} \frac{g^{(n)}(\lambda)}{n!} (T - \lambda \cdot Id)^n \right] P_{\lambda} \quad (9)$$

626 further:

627 We first note, that summing up to $n \geq m_\lambda$ would not make sense, since the nilpotency relation

$$(T - \lambda \cdot Id)^{m_\lambda} \cdot P_\lambda = [(T - \lambda \cdot Id) \cdot P_\lambda]^{m_\lambda} = 0.$$

628 tells us that any such term would be zero.

629 Furthermore, the factor $(T - \lambda \cdot Id)^k$ can be considered to be a "ladder operator" acting on a
 630 basis $\{w_i^\lambda\}_{i=1}^{m_\lambda}$ of generalized eigenvectors spanning the generalized eigenspace associated to the
 631 eigenvalue λ : It acts as

$$(T - \lambda \cdot Id) \cdot P_\lambda \cdot w_i^\lambda = w_{i+1}^\lambda,$$

632 as discussed in Appendix B.

633 The values $g^{(n)}(\lambda)$ in (9) can then be interpreting as weighing the individual "permutations" $(T - \lambda \cdot$
 634 $Id)^n P_\lambda$ of basis elements in the generalized eigenspace associated to the eigenvalue λ .

635 Finally we note that the spectrum of the new operator $g(T)$ is given as

$$\sigma(g(T)) = g(\sigma(T)) := \{g(\lambda) : \lambda \in \sigma(T)\}.$$

636 This is known as the "Spectral Mapping Theorem" [30].

637 **Compatibility with Algebraic relations** Here we prove compatibility of the holomorphic func-
 638 tional calculus with algebraic relations. For ease in writing, we will use the notation

$$(z \cdot Id - T)^{-1} \equiv \frac{1}{z - T}$$

639 below

640 Let us begin with monomials:

641 **Lemma E.1.** Applying the function $g(\lambda) = \lambda^k$ to T yields T^k .

642 *Proof.* We want to prove that

$$\frac{1}{2\pi i} \oint_{\Gamma} \frac{z^k}{z - T} dz = T^k$$

643 To this end, we use the Neumann series characterisation of the resolvent [52]

$$(z - T)^{-1} = \frac{1}{z} \sum_{n=0}^{\infty} \left(\frac{T}{z}\right)^n,$$

644 which is valid for $|z| > \|T\|$. Substituting with $g(\lambda) = \lambda^k$ yields

$$g(T) = \frac{1}{2\pi i} \oint_{\Gamma} \left(\sum_{n=0}^{\infty} \frac{T^n}{z^{n+1-k}} \right) dz$$

645 which we may rewrite as

$$\sum_{n=0}^{\infty} \left(\frac{1}{2\pi i} \oint_{\Gamma} \frac{dz}{z^{n+1-k}} \right) = \sum_{n=0}^{\infty} T^n \cdot \delta_{nk} = T^k$$

646 Here we used the relation (c.f. e.g. [2])

$$\frac{1}{2\pi i} \oint_{\Gamma} \frac{dz}{z^{n+1-k}} = \delta_{nk}.$$

647

□

648 Next we prove that the holomorphic functional calculus is also consistent with inversion:

649 **Lemma E.2.** If y is not an eigenvalue of T , applying the function $g(\lambda) = \left(\frac{1}{\lambda-y}\right)^k$ to T yields

$$g(T) = [(T - y \cdot Id)^{-1}]^k.$$

650 For ease in notation we will write $[(T - y \cdot Id)^{-1}]^k \equiv (T - y \cdot Id)^{-k}$.

651 *Proof.* What want to prove, is thus the equality

$$(y \cdot Id - T)^{-k} := \frac{1}{2\pi i} \oint_{\Gamma} (y - z)^{-k} \cdot (zId - T)^{-1} dz,$$

652 We first note that for the resolvent $R_y(T) = (T - y \cdot Id)^{-1}$ we may write

$$R_x(T) = \sum_{n=0}^{\infty} (x - y)^n (-1)^n R_y(T)^{n+1}$$

653 for $|x - y| \leq \|R_y(T)\|$ using standard results in matrix analysis (namely the 'Neumann Characterisation of the Resolvent' which is obtained by repeated application of a resolvent identity; c.f. [46] or [52] for more details). We thus find

$$\frac{1}{2\pi i} \oint_{\Gamma} \left(\frac{1}{y-z}\right)^k \frac{1}{zId - T} dz = \frac{1}{2\pi i} \oint_{\Gamma} \left(\frac{1}{y-z}\right)^k \sum_{n=0}^{\infty} (y-z)^n R_y(T)^{n+1}.$$

656 Using the fact that

$$\frac{1}{2\pi i} \oint_{\Gamma} (z - y)^{n-k-1} dz = \delta_{nk}$$

657 then yields the claim. □

658 **F Proof of Theorem 3.2**

659 To increase readability, we here use the notation

$$\Delta \equiv L^{\text{in}}.$$

660 In this section, we then prove Theorem 3.2. For convenience, we first restate the result – together with the definitions leading up to it – again:

662 **Definition F.1.** Denote by $\underline{\mathcal{G}}$ the set of reaches in G_{high} . We give this set a graph structure as follows:
663 Let R and P be elements of $\underline{\mathcal{G}}$ (i.e. reaches in G_{high}). We define the real number

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

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

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

665 and assign \underline{W}_{RP} as weight to such edges. Node weights of limit nodes are defined similarly as aggregated weights of all nodes r (in G) contained in the reach R as

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

667 In order to translate signals between the original graph G and the limit description $\underline{\mathcal{G}}$, we need translation operators mapping signals from one graph to the other:
668

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

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

672 The upsampling operator J^\uparrow is defined as

$$J^\uparrow u = \sum_R u_R \cdot \mathbb{1}_R; \quad (10)$$

673 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 reaches
674 in G_{high} .

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

676 **Theorem F.3.** We have $R_z(\Delta) \rightarrow J^\uparrow R_z(\underline{\Delta}) J^\downarrow$ as the weight scale c increases. Explicitly,

$$\|R_z(\Delta) - J^\uparrow R_z(\underline{\Delta}) J^\downarrow\| \rightarrow 0 \text{ as } c \rightarrow \infty$$

677 holds.

678 The proof closely follows that of the corresponding result in [33].

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

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

680 the resolvents corresponding to Δ , Δ_{high} and Δ_{regular} respectively.
681 Our first goal is establishing that we may write

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

682 This will follow as a consequence of what is called the second resolvent formula [52]:

683 "Given operators A, B , we may write

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

684 In our case, this translates to

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

685 or equivalently

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

686 Multiplying with $[Id + R_z(\Delta_{\text{high}})\Delta_{\text{regular}}]^{-1}$ from the left then yields

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

687 as desired.

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

689

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

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

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

$$(\Delta_{\text{high}} + \Delta_{\text{regular}} - zId)x = 0$$

692 which is precisely to say that

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

693 But since Δ is a graph Laplacian, it only has eigenvalues with non-negative real part [54]. Hence we
 694 have reached our contradiction and established

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

695

696 Our next step is to establish that

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

697 where P_0^{high} is the spectral projection onto the eigenspace corresponding to the lowest lying eigenvalue
 698 $\lambda_0(\Delta_{high}) = 0$ of Δ_{high} .

699 Indeed, using the spectral characterization of the holomorphic functional calculus, we may write

$$g(\Delta_{high}) = \sum_{\lambda \in \sigma(\Delta_{high})} g(\lambda)P_\lambda + \sum_{\lambda \in \sigma(\Delta_{high})} \left[\sum_{n=1}^{m_\lambda-1} \frac{g^{(n)}(\lambda)}{n!} (\Delta_{high} - \lambda \cdot Id)^n \right] P_\lambda.$$

700 with

$$g(\lambda) = \frac{1}{\lambda - z}.$$

701 Scaling the operator Δ_{high} as

$$\Delta_{high} \mapsto c \cdot \Delta_{high}$$

702 also scales all corresponding eigenvalues λ as $\lambda \mapsto c \cdot \lambda$, while leaving the spectral projections P_λ
 703 invariant [30]. Thus taking the limit $c \rightarrow \infty$ we indeed find

$$\lim_{c \rightarrow \infty} g(c \cdot \Delta_{high}) = \frac{P_0^{high}}{-z}.$$

704 Our next task is to use this result in order to show that the difference

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

705 goes to zero as $c \rightarrow \infty$.

706 To this end we first note that the relation

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

707 provided to us by the second resolvent formula, implies

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

708 Thus we have

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

709 With this, we have

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

710 Hence it remains to bound the left hand summand. For this we use the following fact (c.f. [27],
711 Section 5.8. "Condition numbers: inverses and linear systems"):

712

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

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

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

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

715 For c sufficiently large, we have

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

716 so that we may estimate

$$\begin{aligned}
& \left\| \left[Id + \Delta_{regular} \frac{P_0^{high}}{-z} \right]^{-1} - [Id + \Delta_{regular} R_z(\Delta_{high})]^{-1} \right\| \\
&\leq 2 \cdot (1 + \|\Delta_{regular}\|) \cdot \left\| \frac{P_0^{high}}{-z} - R_z(\Delta_{high}) \right\| \\
&\rightarrow 0
\end{aligned}$$

717 Thus we have now established

$$\left| \left[Id + \frac{P_0^{high}}{-z} \Delta_{regular} \right]^{-1} \cdot \frac{P_0^{high}}{-z} - R_z(\Delta) \right| \rightarrow 0.$$

718

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

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

720 with $J^\uparrow, \underline{\Delta}, J^\downarrow$ as defined above. To this end, we first note that since the left-kernel and right-kernel
 721 of Δ_{high} are the same (since in-degrees are the same as out degrees), we have

$$J^\uparrow \cdot J^\downarrow = P_0^{\text{high}} \quad (11)$$

722 and

$$J^\downarrow \cdot J^\uparrow = Id_{\underline{G}}. \quad (12)$$

723 Indeed, the relation (11) follows from the fact that the eigenspace corresponding to the eigenvalue zero
 724 is spanned by the vectors $\{\mathbb{1}_R\}_R$, with $\{R\}$ the reaches of G_{high} (c.f. [54]). Equation (12) follows
 725 from the fact that

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

726 With this we have

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

727 To proceed, set

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

728 and

$$\mathcal{X} = \left[P_0^{\text{high}} \Delta_{\text{regular}} - zId \right]^{-1} P_0^{\text{high}} x.$$

729 Then

$$\left[P_0^{\text{high}} \Delta_{\text{regular}} - zId \right] \mathcal{X} = P_0^{\text{high}} x$$

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

$$J^\uparrow J^\downarrow (\Delta_{\text{regular}} - zId) J^\uparrow J^\downarrow \mathcal{X} = J^\uparrow J^\downarrow x.$$

731 Multiplying with J^\downarrow from the left yields

$$J^\downarrow (\Delta_{\text{regular}} - zId) J^\uparrow J^\downarrow \mathcal{X} = J^\downarrow x.$$

732 Thus we have

$$(J^\downarrow \Delta_{\text{regular}} J^\uparrow - zId) J^\uparrow J^\downarrow \mathcal{X} = J^\downarrow x.$$

733 This – in turn – implies

$$J^\uparrow J^\downarrow \mathcal{X} = \left[J^\downarrow \Delta_{\text{regular}} J^\uparrow - zId \right]^{-1} J^\downarrow x.$$

734 Using

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

735 we then have

$$\mathcal{X} = J^\uparrow \left[J^\downarrow \Delta_{\text{regular}} J^\uparrow - zId \right]^{-1} J^\downarrow x.$$

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

738 As a corollary, we find

739 **Corollary F.4.** We have

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

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

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

741 \square

742 This thus establishes Theorem 3.3.

743 **G Stability under Scale Variations**

744 Here we provide details on the scale-invariance results discussed in Section 4; most notably Theorem
745 4.2.

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

748 **Lemma G.1.** Let $\underline{R}_z := (\underline{\Delta} - zId)^{-1}$ and $R_z := (\Delta - zId)^{-1}$. For any natural number k , we have
749

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

750 for

$$\|R_z(\Delta)\|, \|R_z(\underline{\Delta})\| \leq A$$

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

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

752 Iterating this, using the fact that $\|R_z(\Delta)\|$ stays bounded as $c \rightarrow \infty$, since

$$\|R_z(\Delta)\| \rightarrow \|J^\uparrow R_z(\underline{\Delta}) J^\downarrow\| \leq A$$

753 for some constant A together with $\|J^\uparrow\|, \|J^\downarrow\| \leq 1$ and

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

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

755 □

756 Hence let us now prove a node-level stability result:

757 **Theorem G.2.** Let Φ_L and $\underline{\Phi}_L$ be the maps associated to Dir-ResolvNets with the same learned
758 weight matrices and biases but deployed on graphs G and \underline{G} as defined in Section 3.3.2. We have

$$\|\Phi_L(X) - J^\uparrow \underline{\Phi}_L(J^\downarrow X)\|_2 \leq (C_1(\mathcal{W}, A) \cdot \|X\|_2 + C_2(\mathcal{W}, \mathcal{B}, A)) \cdot \|R_z(\Delta) - J^\uparrow R_z(\underline{\Delta}) J^\downarrow\| \quad (13)$$

759 with A a constant such that

$$\|R_z(\Delta)\|, \|R_z(\underline{\Delta})\| \leq A.$$

760 *Proof.* Let us define

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

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

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

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

764 and

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

765 Here, since bias terms are proportional to constant vectors on the graphs, we have

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

766 and

$$J^\uparrow \underline{B} = B \quad (14)$$

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

768 We then have

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

769 Here we used the fact that since $\rho(\cdot)$ maps positive entries to positive entries and acts pointwise, it
770 commutes with J^\uparrow . We also made use of (14).

771 Using the fact that $\rho(\cdot)$ is Lipschitz-continuous with Lipschitz constant $D = 1$, we can establish

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

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

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

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

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

774 Applying Lemma G.1 yields

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

775 Similarly, one may establish that we have

$$\|X^L\| \leq C(A) \cdot \left(\|B^L\| + \sum_{m=0}^L \left(\prod_{j=0}^m \|\mathscr{W}^{L-1-k}\|_z \right) \|B^{L-1-k}\| + \left(\prod_{\ell=1}^L \|\mathscr{W}^\ell\|_z \right) \cdot \|X\| \right). \quad (15)$$

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

778 which tends to zero.

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

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

781 The last step is only possible because we let the sums over powers of resolvents start at $a = 1$ as
782 opposed to $a = 0$. In the latter case, there would have remained a term $\|X - J^\uparrow J^\downarrow X\|$, which would
783 not decay as $c \rightarrow \infty$.

784 Aggregating terms, we build up the polynomial stability constants of (13) layer by layer, and complete
785 the proof. □

786

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

788 **Theorem G.3.** Denote by Ω the aggregation method introduced in Section 4. With $\mu(G) = \sum_{i=1}^N \mu_i$
789 the total weight of the graph G , we have in the setting of Theorem G.2, that

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

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

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

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

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

793 To this end, we first note that

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

794 Indeed, this follows from the fact that given a reach R in G_{high} , the map J^\uparrow assigns the same feature
795 vector to each node $r \in R \subseteq G$ (c.f. (10)), together with the fact that

$$\mu_R = \sum_{r \in R} \mu_r.$$

796 Thus we have

$$\|\Omega(\Phi_L(X)) - \Omega(\Phi_L(J^\downarrow X))\|_2 = \|\Omega(\Phi_L(X)) - \Omega(J^\uparrow \Phi_L(J^\downarrow X))\|_2.$$

797 Next let us simplify notation and write

$$\mathcal{A} = \Phi_L(X)$$

798 and

$$\mathcal{B} = J^\uparrow \Phi_L(J^\downarrow X)$$

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

$$\|\Omega(\Phi_L(X)) - \Omega(J^\uparrow \Phi_L(J^\downarrow X))\|_2^2 = \sum_{j=1}^F \left(\sum_{i=1}^N (|\mathcal{A}_{ij}| - |\mathcal{B}_{ij}|) \cdot \mu_i \right)^2.$$

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

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

801 Since we have

$$\|\Phi_L(X) - J^\uparrow \Phi_L(J^\downarrow X)\|_2^2 = \sum_{j=1}^F \left(\sum_{i=1}^N |\mathcal{A}_{ij} - \mathcal{B}_{ij}|^2 \cdot \mu_i \right),$$

802 the claim is established. \square

803 H Proof of Theorem 4.1

804 Here we prove Theorem 4.1, which we restate here for convenience:

805 **Theorem H.1.** Suppose for filter banks $\{\Psi_i^{\text{fwd}/\text{fwd}}\}_{I^{\text{fwd}/\text{fwd}}}$ that the matrices $\Psi_i^{\text{fwd}}(T), \Psi_i^{\text{bwd}}(T^*)$ con-
 806 tain only real entries. Then any HoloNet with layer-widths $\{F_\ell\}$ with complex weights & biases and
 807 a non linearity that acts on complex numbers componentwise as $\rho(a + ib) = \tilde{\rho}(a) + i\tilde{\rho}(b)$ can be
 808 exactly represented by a HoloNet of widths $\{2 \cdot F_\ell\}$ utilizing $\tilde{\rho}$ and employing only real weights &
 809 biases.

810 *Proof.* It suffices to prove, that in this setting the update rule

$$X^\ell = \rho \left(\alpha \sum_{i \in I} \Psi_i^{\text{fwd}}(T) \cdot X^{\ell-1} \cdot W_i^{\text{fwd}, \ell} + (1 - \alpha) \sum_{i \in I} \Psi_i^{\text{bwd}}(T^*) \cdot X^{\ell-1} \cdot W_i^{\text{bwd}, \ell} + B^\ell \right).$$

811 can be replaced by purely real weights and biases. For simplicity in notation, let us assume $\alpha = 1$;
 812 the general case follows analogously but with more cluttered notation.

813 Let us write $X = X_{\text{real}} + iX_{\text{imag}}, W = W_{\text{real}} + iW_{\text{imag}}, B = B_{\text{real}} + iB_{\text{imag}}$. Then we have
 814 with $\{\Psi_i^{\text{fwd}}(T)\}_{i \in I}$ purely real, that

$$\begin{aligned} X^\ell &= \rho \left(\sum_{i \in I} \Psi_i^{\text{fwd}}(T) \cdot X^{\ell-1} \cdot W_i^{\text{fwd}, \ell} + B^\ell \right) \\ &= \rho \left(\sum_{i \in I} \Psi_i^{\text{fwd}}(T) \cdot (X_{\text{real}}^{\ell-1} W_{\text{real}, i}^{\text{fwd}, \ell} - X_{\text{imag}}^{\ell-1} W_{\text{imag}, i}^{\text{fwd}, \ell}) + B_{\text{real}}^\ell \right. \\ &\quad \left. + i \left[\sum_{i \in I} \Psi_i^{\text{fwd}}(T) \cdot (X_{\text{real}}^{\ell-1} W_{\text{imag}, i}^{\text{fwd}, \ell} + X_{\text{imag}}^{\ell-1} W_{\text{real}, i}^{\text{fwd}, \ell}) + B_{\text{imag}}^\ell \right] \right) \\ &= \tilde{\rho} \left(\sum_{i \in I} \Psi_i^{\text{fwd}}(T) \cdot (X_{\text{real}}^{\ell-1} W_{\text{real}, i}^{\text{fwd}, \ell} - X_{\text{imag}}^{\ell-1} W_{\text{imag}, i}^{\text{fwd}, \ell}) + B_{\text{real}}^\ell \right) \\ &\quad + i \tilde{\rho} \left(\sum_{i \in I} \Psi_i^{\text{fwd}}(T) \cdot (X_{\text{real}}^{\ell-1} W_{\text{imag}, i}^{\text{fwd}, \ell} + X_{\text{imag}}^{\ell-1} W_{\text{real}, i}^{\text{fwd}, \ell}) + B_{\text{imag}}^\ell \right) \end{aligned}$$

815 The result then immediately follows after using the canonical isomorphism between \mathbb{C}^d and \mathbb{R}^{2d} as

$$X^\ell \cong \begin{pmatrix} X_{\text{real}}^\ell \\ X_{\text{imag}}^\ell \end{pmatrix} = \tilde{\rho} \left[\begin{pmatrix} \sum_{i \in I} \Psi_i^{\text{fwd}}(T) \cdot (X_{\text{real}}^{\ell-1} W_{\text{real}, i}^{\text{fwd}, \ell} - X_{\text{imag}}^{\ell-1} W_{\text{imag}, i}^{\text{fwd}, \ell}) + B_{\text{real}}^\ell \\ \sum_{i \in I} \Psi_i^{\text{fwd}}(T) \cdot (X_{\text{real}}^{\ell-1} W_{\text{imag}, i}^{\text{fwd}, \ell} + X_{\text{imag}}^{\ell-1} W_{\text{real}, i}^{\text{fwd}, \ell}) + B_{\text{imag}}^\ell \end{pmatrix} \right]. \quad (16)$$

816 The above layer update

$$\begin{pmatrix} X_{real}^{\ell-1} \\ X_{imag}^{\ell-1} \end{pmatrix} \xrightarrow{(16)} \begin{pmatrix} X_{real}^{\ell} \\ X_{imag}^{\ell} \end{pmatrix}$$

817 can then clearly be realised by a real network as described in Theorem 4.1.

818

□

819 I Additional Details on Experiments:

820 I.1 FaberNet: Node Classification

821 **Datasets:** We evaluate on the task of node classification on several directed benchmark datasets with
 822 high homophily: Chameleon & Squirrel [44], Arxiv-Year [29], Snap-Patents [37] and Roman-Empire
 823 [45]. These datasets are highly heterophilic (edge homophily smaller than 0.25).

Table 5: Node Classification Datasets: Statistics

DATASET	# NODES	# EDGES	# FEAT.	# C	UNID. EDGES	EDGE HOM.
CHAMELEON	2,277	36,101	2,325	5	85.01%	0.235
SQUIRREL	5,201	217,073	2,089	5	90.60%	0.223
ARXIV-YEAR	169,343	1,166,243	128	40	99.27%	0.221
SNAP-PATENTS	2,923,922	13,975,791	269	5	99.98%	0.218
ROMAN-EMPIRE	22,662	44,363	300	18	65.24%	0.050

824 **Experimental Setup** All experiments are conducted on a machine with NVIDIA A4000 GPU with
 825 16GB of memory, safe for experiments on snap-patents which have been performed on a machine
 826 with one NVIDIA Quadro RTX 8000 with 48GB of memory. We closely follow the experimental
 827 setup of [47]. In all experiments, we use the Adam optimizer and train the model for 10000 epochs,
 828 using early stopping on the validation accuracy with a patience of 200 for all datasets apart from
 829 Chameleon and Squirrel, for which we use a patience of 400. For OGBN-Arxiv we use the fixed split
 830 provided by OGB [29], for Chameleon and Squirrel we use the fixed GEOM-GCN splits [44], for
 831 Arxiv-Year and Snap-Patents we use the splits provided in [37], while for Roman-Empire we use the
 832 splits from [45].

833 **Baselines Results:** Results for MLP, GCN, H₂GCN, GPR-GNN and LINKX were taken from
 834 (author?). Results for Gradient Gating are taken from their paper [51]. Results for GloGNN
 835 are taken from their paper [36]. Results on Roman-Empire are taken from [45] for GCN, H₂GCN,
 836 GPR-GNN, FSGNN and GloGNN and from [47] for MLP, LINKX, ACM-GCN and Gradient Gating.
 837 Results for FSGNN are taken from [40] for Actor, Squirrel and Chameleon, and from [47] for results
 838 on Arxiv-year and Snap-Patents. Results for DiGCN, MagNet are taken from [47]. Results for
 839 DirGNN were obtained via a re-implementation; using the official codebase and hyperparameters
 840 specified in [47]. Note that – as detailed in [47] – the reported results for DirGNN correspond to a
 841 best-of-three report over directed version of GCN [31], GAT [56] and Sage [21].

842 **Hyperparameters:** Following the setup of [47], our search space for generic hyperparameters
 843 is given by varying the learning rate $lr \in \{0.01, 0.005, 0.001, , 0.0005\}$, the hidden dimension
 844 over $F \in \{32, 64, 128, 256, 512\}$, ne number of layers over $L \in \{2, 3, 4, 5, 6\}$, jumping knowledge
 845 connections over $jk \in \{max, cat, none\}$ layer-wise normalization in $norm \in \{True, False\}$,
 846 patience as $patience \in \{200, 400\}$ and dropout as $p \in \{0, 0.2, 0.4, 0.6, 0.8, 1\}$.

847 In practice we take the parameters of Table 6 as frozen and given by [47]. We then optimize over the
 848 custom hyperparameters pertaining to our method. To this end, we vary the maximal order of our
 849 Faber polynomials $\{\Psi_i\}_{i=0,1}^K$ as $K \in \{1, 2, 3, 4, 5\}$. Note that we also discount higher order terms
 850 with a regularization $\sim 1/2^i$, as this improved results experimentally. We thus have $\Psi_k = \lambda^k/2^k$.
 851 The type of weights & biases is varied over $parameters \in \{\mathbb{R}, \mathbb{C}\}$. The non-linearity is varied over
 852 $\{\cdot|_{\mathbb{C}}, \text{ReLU}\}$, with $|a + ib|_{\mathbb{C}} = |a| + i|b|$. The parameter α is varied as $\alpha \in \{0, 0.5, 1\}$ as in [47].
 853 The zero-order Faber polynomial $\Psi_0(\lambda) = 1$ is either included or discarded, as discussed in Section 4

Table 6: Selected Generic Hyperparameters

DATASET	lr	L	PATIENCE	F	NORM	P	JK
CHAMELEON	0.005	5	400	128	TRUE	0	CAT
SQUIRREL	0.01	4	400	128	TRUE	0	MAX
ARXIV-YEAR	0.005	6	200	256	FALSE	0	CAT
SNAP-PATENTS	0.01	5	200	32	TRUE	0	MAX
ROMAN-EMPIRE	0.01	5	200	256	FALSE	0.2	CAT

854 and weight decay parameters for real- and imaginary weights are varied over $\lambda_{\text{real}}, \lambda_{\text{imag}} \in \{0, 0.1, 1\}$.
 855 Final selected hyperparameters are listed in Table 7.

Table 7: Selected Custom Hyperparameters

DATASET	K	PARAMETERS	NON.-LIN.	α	Ψ_0	λ_{REAL}	λ_{IMAG}
CHAMELEON	4	\mathbb{C}	$ \cdot _{\mathbb{C}}$	0	NO	1	0
SQUIRREL	5	\mathbb{C}	$ \cdot _{\mathbb{C}}$	0	NO	0.1	0.1
ARXIV-YEAR	2	\mathbb{R}	RELU	0.5	NO	0.1	N.A.
SNAP-PATENTS	2	\mathbb{R}	RELU	0.5	NO	0.1	N.A.
ROMAN-EMPIRE	1	\mathbb{R}	RELU	0.5	YES	0.1	N.A.

856 I.2 Dir-ResolvNet: Digraph Regression and Scale Insensitivity

857 **Dataset:** The dataset we consider is the **QM7** dataset, introduced in [7, 50]. This dataset contains
 858 descriptions of 7165 organic molecules, each with up to seven heavy atoms, with all non-hydrogen
 859 atoms being considered heavy. In the dataset, a given molecule is represented by its Coulomb matrix
 860 W , whose off-diagonal elements

$$W_{ij} = \frac{Z_i Z_j}{|\vec{x}_i - \vec{x}_j|} \quad (17)$$

861 correspond to the Coulomb-repulsion between atoms i and j . We discard diagonal entries of Coulomb
 862 matrices; which would encode a polynomial fit of atomic energies to nuclear charge [50].

863 To each molecule an atomization energy - calculated via density functional theory - is associated. The
 864 objective is to predict this quantity. The performance metric is mean absolute error. Numerically,
 865 atomization energies are negative numbers in the range -600 to -2200 . The associated unit is
 866 $[kcal/mol]$.

867 For each atom in any given molecular graph, the individual Cartesian coordinates \vec{x}_i and the atomic
 868 charge Z_i are also accessible individually.

869 In order to induce directedness, we modify the Coulomb weights (17): Weights *only* from heavy
 870 atoms to *atoms outside this heavy atom’s respective immediate hydrogen cloud* are modified as

$$W_{ij} := \frac{Z_i^{\text{outside}} \cdot (Z_j^{\text{heavy}} - 1)}{|\vec{x}_i - \vec{x}_j|}. \quad (18)$$

871 The immediate hydrogen cloud of a given heavy atom, we take to encompass precisely those hydrogen
 872 atoms for which this heavy atom is the closest among all other heavy atoms.

873 This specific choice (18) is made in preparation for the scale insensitivity experiments: The theory
 874 developed in Section 3.3.2 applies to those graphs, where strongly connected subgraphs contain only
 875 nodes for which the in-degree equals the out-degree (where only strong weights are considered when
 876 calculating the respective degrees). The choice (18) facilitates this, as hydrogen-hydrogen weights
 877 and weights between hydrogen and respective closest heavy atom remain symmetric.

878 **Experimental Setup:** We shuffle the dataset and randomly select 1500 molecules for testing. We
 879 then train on the remaining graphs. We run experiments for 5 different random seeds and
 880 report mean and standard deviation.

881 All considered convolutional layers (i.e. for Dir-ResolvNet and baselines) are incorporated into a
 882 two layer deep and fully connected graph convolutional architecture. In each hidden layer, we set the
 883 width (i.e. the hidden feature dimension) to

$$F_1 = F_2 = 64.$$

884 For all baselines, the standard mean-aggregation scheme is employed after the graph-convolutional
 885 layers to generate graph level features. Finally, predictions are generated via an MLP.

886 For Dir-ResolvNet, we take $\alpha = 1$, use real weights and biases and set $z = -1$. These choices are
 887 made for simplicity. Resolvents are thus given as

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

888 As aggregation for our model, we employ the graph level feature aggregation scheme Ω introduced
 889 before Theorem 4.2 in Section 4. Node weights set to atomic charges of individual atoms. Predictions
 890 are then generated via a final MLP with the same specifications as the one used for baselines.

891 **Scale Insensitivity** We then modify (all) molecular graphs in QM7 by deflecting hydrogen atoms
 892 (H) out of their equilibrium positions towards the respective nearest heavy atom. This is possible
 893 since the QM7 dataset also contains the Cartesian coordinates of individual atoms.

894 This introduces a two-scale setting precisely as discussed in section 3.3.2: Edge weights between
 895 heavy atoms remain the same, while Coulomb repulsions between H-atoms and respective nearest
 896 heavy atom increasingly diverge; as is evident from (17).

897 Given an original molecular graph G with node weights $\mu_i = Z_i$, the corresponding limit graph
 898 \underline{G} corresponds to a coarse grained description, where heavy atoms and surrounding H-atoms are
 899 aggregated into single super-nodes in the sense of Section 3.3.2.

900 Mathematically, \underline{G} is obtained by removing all nodes corresponding to H-atoms from G , while adding
 901 the corresponding charges $Z_H = 1$ to the node-weights of the respective nearest heavy atom. Charges
 902 in (18) are modified similarly to generate the weight matrix \underline{W} .

903 On original molecular graphs, atomic charges are provided via one-hot encodings. For the graph of
 904 methane – consisting of one carbon atom with charge $Z_C = 6$ and four hydrogen atoms of charges
 905 $Z_H = 1$ – the corresponding node-feature-matrix is e.g. given as

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

906 with the non-zero entry in the first row being in the 6th column, in order to encode the charge $Z_C = 6$
 907 for carbon.

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

910 As discussed in Section 3.3.2, node feature matrices are translated as $\underline{X} = J^\downarrow X$. Applying J^\downarrow to
 911 one-hot encoded atomic charges yields (normalized) bag-of-word embeddings on \underline{G} : Individual
 912 entries of feature vectors encode how much of the total charge of the super-node is contributed by
 913 individual atom-types. In the example of methane, the limit graph \underline{G} consists of a single node with
 914 node-weight

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

915 The feature matrix

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

916 is a single row-vector given as

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