

ON MULTI-SCALE GRAPH REPRESENTATION LEARNING

Anonymous authors

Paper under double-blind review

1 INTRODUCTION

Graph neural networks (GNNs) play a fundamental role in modern computational biology, where they often form the backbone for both subcellular tasks such as protein structure prediction (Jumper et al., 2021) as well as cell type annotation (Shao et al., 2021) and modeling of protein protein interactions (Liu et al., 2019) at the multi-cellular and organism wide levels. An underexplored drawback of common GNN methods, however, is that they are not inherently multiscale consistent: Two graphs describing the same object or situation at different resolution scales are assigned vastly different latent representations. This prevents graph networks from generating data representations that are consistent across scales. It also complicates the integration of representations at the molecular scale with those generated at the biological scale. Here we discuss why existing GNNs struggle with multiscale consistency and show how to overcome this problem by modifying the message passing paradigm within GNNs.

2 STANDARD GNNs ARE NOT MULTI-SCALE CONSISTENT

To illustrate that standard GNNs are unable to consistently integrate multiple scales, we utilize the QM7 dataset (Rupp et al., 2012). Here, molecular atomization energies of organic molecules (containing both hydrogen and heavy atoms) are to be predicted. Each molecule is represented by an adjacency matrix with entries $A_{ij} = Z_i Z_j |\vec{x}_i - \vec{x}_j|^{-1}$ given as Coulomb energies of atoms i, j .

From a physical perspective, describing a molecule at the level of interacting atoms corresponds to a specific choice of resolution scale, where interactions of individual protons and neutrons inside individual atoms are discarded. To test the multi-scale consistency of GNNs we additionally also consider a version of QM7 where we further lower the resolution scale: Here we aggregate each heavy atomic core together with its surrounding (single-proton) hydrogen atoms into super-nodes.

To showcase the failure of GNNs to consistently incorporate multiple scales, we confront models during inference with a version of QM7 on a scale different from the one they were trained on. As Table 1 details, mean-absolute-errors (MAEs) increase significantly when going from a same-resolution setting to a cross-resolution setting. None of the considered standard architectures (including multi-scale methods (SAG-M – PushNet)) consistently handles multiple scales. We can trace this back to the latent embeddings F and \underline{F} that are being generated for original- $\{G\}$ and coarsified graphs $\{\underline{G}\}$: For models of Table 1 on average $10 \lesssim \|F - \underline{F}\| \lesssim 10^4$ (c.f. also Fig. 2): Latent representations of graphs describing the same object at different resolutions differ significantly.

To understand this behavior, we interpolate between fine and coarse resolution: Original graphs $\{G\}$ of QM7 are modified ($\{G_\omega\}$) by moving hydrogen atoms towards their corresponding heavy atom by a factor of $\omega \geq 1$ (i.e. $\text{dist}_{\text{new}} = \text{dist}_{\text{equilib.}}/\omega$). For $\omega \rightarrow \infty$, they arrive at the respective heavy atom ($\{\underline{G}\}$). In Fig. 2, we compare the latent distance between the coarse embeddings \underline{F} and the embeddings for the intermediate graph F_ω . Embeddings F_ω do *not* converge to the coarse embeddings \underline{F} . Since the convergence of graph-sequence G_ω to the limit graph \underline{G} is not turned into a convergence of latent embeddings $F_\omega \rightarrow \underline{F}$ we conclude: **GNNs are not continuous**. This discontinuity explains why GNNs can map similar graphs (describing the same object at different resolutions) to different latent representations.

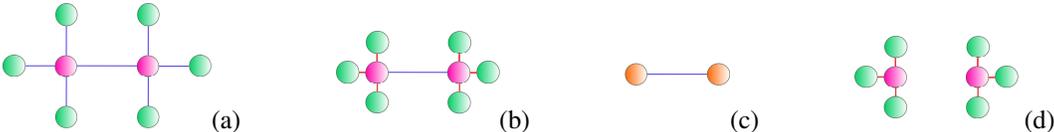
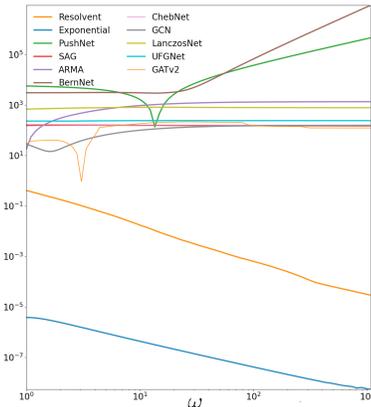


Figure 1: (a) G of QM7 (b) Modified G_ω (c) Coarsified \underline{G} (d) Effective propagation in GCN

Table 1: Regression using high- and low-resolution QM7

Mean Absolute Error (\downarrow) on QM7 [kcal/mol]				
Training	High Resolution		Low Resolution	
	Low Resolution	High Resolution	Low Resolution	High Resolution
Inference				
GCN	125.34 \pm 2.47	63.17 \pm 0.92	67.75 \pm 3.73	380.51 \pm 30.33
GATv2	415.09 \pm 96.57	48.41 \pm 19.20	60.01 \pm 3.34	245.03 \pm 90.97
ChebNet	568.47 \pm 37.70	64.63 \pm 1.21	64.90 \pm 4.55	339.64 \pm 101.30
SAG	542.16 \pm 27.33	68.43 \pm 1.93	104.20 \pm 3.92	506.75 \pm 60.57
BernNet	765.22 \pm 495.28	83.76 \pm 21.75	90.52 \pm 37.17	594.62 \pm 341.55
SAG-M	285.53 \pm 95.54	66.22 \pm 4.51	73.57 \pm 14.57	307.67 \pm 77.24
UFGNet	620.21 \pm 4.80	13.71 \pm 1.05	24.53 \pm 4.80	156.44 \pm 156.44
Lanczos	939.87 \pm 16.35	10.55 \pm 3.22	83.11 \pm 5.27	654.61 \pm 529.13
PushNet	2442.59 \pm 303.27	60.94 \pm 1.83	69.25 \pm 3.11	124.08 \pm 3.94
Resolvent	16.54 \pm 3.01	16.53 \pm 3.03	15.79 \pm 0.98	13.80 \pm 1.34
Exponential	16.37 \pm 1.71	16.36 \pm 2.16	16.25 \pm 1.41	16.25 \pm 1.41

Figure 2: Latent distance $\|F_\omega - \underline{F}\|$

To understand this discontinuity, we exemplarily investigate (GCNs) (Kipf & Welling, 2017). There the layer-wise update acts as $X \mapsto \hat{A}XW$, with the feature matrix $X \in \mathbb{R}^{N \times F}$ (N nodes; latent dimension F), the weight matrix $W \in \mathbb{R}^{F \times F}$ and the *renormalized* adjacency matrix $\hat{A} \in \mathbb{R}^{N \times N}$. As hydrogen atoms move closer to the heavy atoms, the entries $\hat{A}_{\text{heavy,heavy}}$ in $\hat{A}_{ij} \sim A_{ij}/\sqrt{d_i d_j}$ tend to zero (as degrees d_{heavy} tend to infinity). Thus communication between heavy atoms becomes severely disrupted. Information only propagates along a increasingly disconnected effective graph (Fig.1(d)).

3 GLOBAL LAPLACE PROPAGATION FACILITATES SCALE-CONSISTENCY

To avoid a disconnected effective propagation graph as in Fig. 1 (d) we modify the message passing paradigm in GCN: To connect the information flows over G_ω and \underline{G} we observe that features in G_ω should equalize faster between nodes connected by large edge weights. When such a large weight tends to infinity, features between strongly connected nodes are then equalized immediately, so that entire strongly connected clusters exactly behave as the single nodes in \underline{G} .

Noting that this is exactly the behavior that heat dissipating over a graph exhibits, we make use of the the heat diffusion equation $dX(t)/dt = -L \cdot X(t)$ (with Graph Laplacian L and time t) and the structure $X(t) = e^{-Lt} \cdot X(0)$ of its solutions, when designing our graph networks:

Definition 3.1. Let $\hat{\psi}$ be a bounded (generalized) function defined on $[0, \infty)$. A **Global Laplacian Propagation Matrix** $\psi(L)$ is any matrix arising as $\psi(L) := \int_0^\infty e^{-tL} \hat{\psi}(t) dt$.

Thus $\psi(L)$ represent a weighted sum of diffusion flows that have progressed to various times. Specifically, if we choose the Dirac distribution $\hat{\psi}_{\delta_{t_k}}(t) := \delta(t - t_k)$ as the weightinf function $\hat{\psi}_k$, we obtain **exponential** matrices $\psi_k(L) = \int_0^\infty \delta(t - t_k) e^{-tL} dt = e^{-t_k L}$ and $\hat{\psi}_k := (-t)^{k-1} e^{-\lambda t}$ to get powers of **resolvents** $\psi_k(L) = [(zId + L)^{-1}]^k$. The propagation matrix is then used instead of the adjacency matrix in each layer of the GNN leading to the update rule $X \mapsto \sum_k \psi_k(L) X W_k$.

As we prove in Appendix C.3, we indeed have $\|F_\omega - \underline{F}\| \rightarrow 0$ as $\omega \rightarrow \infty$, for such networks based on global Laplacian propagation matrices. This behaviour can clearly also be observed for the examples of exponential- and resolvent propagation matrices in Fig. 2. As ω increases, the distance between latent embeddings tends to zero. Thus these networks are indeed continuous.

In the previous Section 2, we had identified the discontinuity of standard GNNs as the obstruction to consistently incorporating multiple scales. This explained their sub-par performance in Table 1. Since networks based on global Laplacian propagation schemes are continuous, we thus expect a consistent incorporation of scales, as well as a good performance in cross resolution setting. This is exactly what we observe in Table 1: MAEs of GNNs based on global Laplacian propagation schemes (using either exponential or resolvent matrices) do not increase when going from a same- to a cross-resolution setting; MAEs of such methods are lower than those of standard graph learning methods by factors of order 10^1 to 10^2 . Hence these methods indeed do consistently incorporate varying scales.

108 MEANINGFULNESS STATEMENT
109

110 A model generating meaningful representations of life should be able to consistently represent all
111 aspects of life across all relevant scales, starting from the molecular level all the way up to the
112 biological level. While graph neural networks have emerged as a popular network architecture
113 for biological problems at any individual scale, we show that they are not directly suitable to
114 facilitate connections *between* respective scales. To remedy this and facilitate progress in eventually
115 continuously traversing between the molecular and the biological scale, we propose a new propagation
116 scheme that allows graph neural networks to indeed incorporate multiple scales.

117
118 REFERENCES

- 119 Wolfgang Arendt. APPROXIMATION OF DEGENERATE SEMIGROUPS. *Taiwanese Journal*
120 *of Mathematics*, 5(2):279 – 295, 2001. doi: 10.11650/twj/1500407337. URL <https://doi.org/10.11650/twj/1500407337>.
- 121
122
- 123 Filippo Maria Bianchi, Daniele Grattarola, Lorenzo Francesco Livi, and Cesare Alippi. Graph neural
124 networks with convolutional arma filters. *IEEE Transactions on Pattern Analysis and Machine*
125 *Intelligence*, 44:3496–3507, 2019.
- 126
127
- 128 L. C. Blum and J.-L. Reymond. 970 million druglike small molecules for virtual screening in the
129 chemical universe database GDB-13. *J. Am. Chem. Soc.*, 131:8732, 2009.
- 130
131
- 132 Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks? In *The Tenth In-*
133 *ternational Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.*
134 *OpenReview.net*, 2022. URL <https://openreview.net/forum?id=F72ximsx7C1>.
- 135
136
- 137 Julian Busch, Jiaying Pi, and Thomas Seidl. Pushnet: Efficient and adaptive neural message passing.
138 In Giuseppe De Giacomo, Alejandro Catalá, Bistra Dilkina, Michela Milano, Senén Barro, Alberto
139 Bugarín, and Jérôme Lang (eds.), *ECAI 2020 - 24th European Conference on Artificial Intelligence,*
140 *29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020*
141 *- Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020),*
142 *volume 325 of Frontiers in Artificial Intelligence and Applications*, pp. 1039–1046. IOS Press,
143 2020. doi: 10.3233/FAIA200199. URL <https://doi.org/10.3233/FAIA200199>.
- 144
145
- 146 Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on
147 graphs with fast localized spectral filtering. *Advances in neural information processing systems*,
148 29, 2016.
- 149
150
- 151 Mingguo He, Zhewei Wei, Zengfeng Huang, and Hongteng Xu. Bernnet: Learning arbitrary
152 graph spectral filters via bernstein approximation. In Marc’Aurelio Ranzato, Alina
153 Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Ad-*
154 *vances in Neural Information Processing Systems 34: Annual Conference on Neural In-*
155 *formation Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.
156 14239–14251, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/76f1cfd7754a6e4fc3281bcccb3d0902-Abstract.html>.
- 157
158
- 159 Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- 160
161
- 162 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
163 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate
164 protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.
- 165
166
- 167 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
168 In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-*
169 *26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=SJU4ayYgl>.
- 170
171
- 172 Christian Koke. Strong connectivity in graphs: Norm resolvent convergence to effective descriptions,
173 2024.

- 162 Christian Koke and Daniel Cremers. Holonets: Spectral convolutions do extend to directed graphs.
 163 In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=EhmEwfavOW>.
 164
 165
- 166 Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. In Kamalika Chaudhuri
 167 and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine*
 168 *Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 3734–3743. PMLR,
 169 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/lee19c.html>.
 170
- 171 Renjie Liao, Zhizhen Zhao, Raquel Urtasun, and Richard S. Zemel. Lanczosnet: Multi-scale deep
 172 graph convolutional networks. In *7th International Conference on Learning Representations,*
 173 *ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=BkedznAqKQ>.
 174
- 175 Zhiyu Liu, Chia-Hsiu Wu, Yujia Zhang, et al. Graph neural network-based prediction of protein-
 176 protein interactions in ppi networks. *Nature Communications*, 10(1):1–10, 2019. doi: 10.1038/
 177 s41467-019-09308-3.
 178
- 179 M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld. Fast and accurate modeling of
 180 molecular atomization energies with machine learning. *Physical Review Letters*, 108:058301,
 181 2012.
 182
- 183 Xin Shao, Haihong Yang, Xiang Zhuang, Jie Liao, Penghui Yang, Junyun Cheng, Xiaoyan Lu, Huajun
 184 Chen, and Xiaohui Fan. scdeepsort: a pre-trained cell-type annotation method for single-cell
 185 transcriptomics using deep learning with a weighted graph neural network. *Nucleic acids research*,
 186 49(21):e122–e122, 2021.
- 187 T. Tao. *An Introduction to Measure Theory*. Graduate studies in mathematics. American Mathematical
 188 Society, 2013. ISBN 9781470409227. URL [https://books.google.de/books?id=](https://books.google.de/books?id=SPGJjwEACAAJ)
 189 [SPGJjwEACAAJ](https://books.google.de/books?id=SPGJjwEACAAJ).
 190
- 191 Gerald Teschl. *Mathematical Methods in Quantum Mechanics*. American Mathematical Society,
 192 2014.
 193
- 194 Xuebin Zheng, Bingxin Zhou, Junbin Gao, Yuguang Wang, Pietro Lió, Ming Li, and Guido Mon-
 195 túfar. How framelets enhance graph neural networks. In Marina Meila and Tong Zhang (eds.),
 196 *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July*
 197 *2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 12761–12771.
 198 PMLR, 2021. URL <http://proceedings.mlr.press/v139/zheng21c.html>.
 199

200 A COARSE-GRAINING GRAPHS

201
 202 In this Appendix we illustrate:
 203

$$204 \quad \|(L + Id)^{-1} - J^\uparrow(\underline{L} + Id)^{-1}J^\downarrow\| \lesssim 1/\lambda_1(\Delta_{\text{high}}).$$

207 which – as we will see – also implies

$$208 \quad \|e^{-tL} - J^\uparrow e^{-t\underline{L}} J^\downarrow\| \lesssim 1/w_{\text{high}}^{\min} \text{ for any } t > 0. \quad (1)$$

209
 210 after noting the linear relation in scaling behaviour $\lambda_1(L_{\text{cluster}}) \sim w_{\text{high}}^{\min}$.
 211
 212

213
 214 This will be used to prove the convergence result of Section 3 in Appendix
 215 For convenience, we restate the definitions leading up to this setting again:

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

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

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

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

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 component R as

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

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:

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

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

The upsampling operator J^\uparrow is defined as

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

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

As proved in (Koke, 2024), we then have the following:

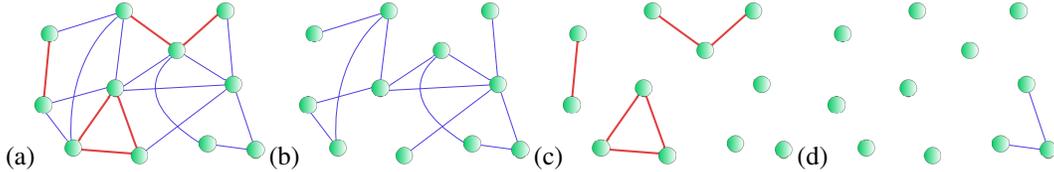


Figure 3: (a) Graph G with $\mathcal{E}_{reg.}$ (blue) & \mathcal{E}_{high} (red); (b) $G_{reg.}$; (c) G_{high} ; (d) $G_{reg., exclusive}$

Theorem A.3. We have

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

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

We here restate the proof for convenience. We use the notation $\Delta = L$.

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

$$R_z(\Delta) = (\Delta - zId)^{-1},$$

$$R_z(\Delta_{high}) = (\Delta_{high} - zId)^{-1}$$

$$R_z(\Delta_{reg.}) = (\Delta_{reg.} - zId)^{-1}$$

the resolvents corresponding to Δ , Δ_{high} and $\Delta_{reg.}$ respectively.

Our first goal is establishing that we may write

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

270 This will follow as a consequence of what is called the second resolvent formula Teschl (2014):

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

$$272 R_z(A + B) - R_z(A) = -R_z(A)BR_z(A + B)."$$

273 In our case, this translates to

$$274 R_z(\Delta) - R_z(\Delta_{high}) = -R_z(\Delta_{high})\Delta_{reg}R_z(\Delta)$$

275 or equivalently

$$276 [Id + R_z(\Delta_{high})\Delta_{reg.}] R_z(\Delta) = R_z(\Delta_{high}).$$

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

$$278 R_z(\Delta) = [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} \cdot R_z(\Delta_{high})$$

279 as desired.

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

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

$$282 [Id + R_z(\Delta_{high})\Delta_{reg.}] x = 0.$$

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

$$284 (\Delta_{high} + \Delta_{reg.} - zId)x = 0$$

285 which is precisely to say that

$$286 (\Delta - zId)x = 0$$

287 But since Δ is a graph Laplacian, it only has non-negative eigenvalues. Hence we have reached our contradiction and established

$$288 R_z(\Delta) = [Id + R_z(\Delta_{high})\Delta_{reg.}]^{-1} R_z(\Delta_{high}).$$

289 Our next step is to establish that

$$290 R_z(\Delta_{high}) \rightarrow \frac{P_0^{high}}{-z},$$

291 where P_0^{high} is the spectral projection onto the eigenspace corresponding to the lowest lying eigenvalue $\lambda_0(\Delta_{high}) = 0$ of Δ_{high} . Indeed, by the spectral theorem for finite dimensional operators (c.f. e.g. Teschl (2014)), we may write

$$292 R_z(\Delta_{high}) \equiv (\Delta_{high} - zId)^{-1} = \sum_{\lambda \in \sigma(\Delta_{high})} \frac{1}{\lambda - z} \cdot P_\lambda^{high}.$$

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

$$294 \left\| R_z(\Delta_{high}) - \frac{P_0^{high}}{-z} \right\| = \left\| \sum_{0 < \lambda \in \sigma(\Delta_{high})} \frac{1}{\lambda - z} \cdot P_\lambda^{high} \right\|;$$

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

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

$$297 \left\| R_z(\Delta_{high}) - \frac{P_0^{high}}{-z} \right\| = \frac{1}{\lambda_1(\Delta_{high}) + |z|}.$$

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

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

$$300 \lambda(S \cdot \Delta) = S \cdot \lambda(\Delta)$$

for any graph Laplacian (in fact any matrix) Δ with eigenvalue λ . Thus we have

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

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

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

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

To this end we first note that the relation

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

provided to us by the second resolvent formula, implies

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

Thus we have

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

With this, we have

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

Hence it remains to bound the left hand summand. For this we use the following fact (c.f. Horn & Johnson (2012), Section 5.8. "Condition numbers: inverses and linear systems"):

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

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

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

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

For S_{high} sufficiently large, we have

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

so that we may estimate

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

Thus we have now established

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

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

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

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

$$J^\uparrow \cdot J^\downarrow = P_0^{\text{high}} \quad (2)$$

and

$$J^\downarrow \cdot J^\uparrow = Id_G. \quad (3)$$

Indeed, the relation (2) follows from the fact that the eigenspace corresponding to the eigenvalue zero is spanned by the vectors $\{\mathbb{1}_R\}_R$, with $\{R\}$ the connected components of G_{high} . Equation (3) follows from the fact that

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

With this we have

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

To proceed, set

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

and

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

Then

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

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

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

Multiplying with J^\downarrow from the left yields

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

Thus we have

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

This – in turn – implies

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

Using

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

we then have

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

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

As a corollary, we find

Corollary A.4. *We have*

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

Proof. This follows directly from the fact that

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

\square

To prove (1), we establish the following theorem:

Theorem A.5. *Consider a graph sequence G_n with $\|(L_n + \lambda Id)^{-1} - \tilde{J}_n(\tilde{L} + \lambda Id)^{-1} J_n\| \rightarrow 0$. Then we have $\|\psi(L_n) - \tilde{J}_n \psi(\tilde{L}) J_n\| \rightarrow 0$ if ψ is complex differentiable and $\lim_{r \rightarrow \infty} \psi(r) = 0$.*

Proof. We make use of the holomorphic functional calculus (c.f. e.g. (Koke & Cremers, 2024)) to establish

$$\|\psi(L) - \tilde{J} \psi(\tilde{L}) J\| \leq \frac{1}{2\pi} \oint_{\Gamma} |\psi(z)| \cdot \|(L - zId)^{-1} - \tilde{J}(\tilde{L} - zId)^{-1} J\| |dz|.$$

Since $\|(L_n + \lambda Id)^{-1} - \tilde{J}_n(\tilde{L} + \lambda Id)^{-1} J_n\| \rightarrow 0$ implies $\|(L_n - zId)^{-1} - \tilde{J}_n(\tilde{L} - zId)^{-1} J_n\| \rightarrow 0$ uniformly (in z) on compact sets (c.f. e.g. Arendt (2001)), we can apply dominated convergence, if we find an majorizing function that is integrable on Γ . But this is ensured by the decay of ψ . \square

Choosing the function ψ to be given as $\psi(z) = e^{-tz}$ then establishes (1).

B GLOBAL LAPLACIAN PROPAGATION MATRICES, GENERALIZED FUNCTIONS, MEASURES AND ALL THAT

In this section we discuss global Laplacian propagation matrices, generalized functions and measures

B.1 COMPLEX MEASURES ON $\mathbb{R}_{\geq 0}$ AND THEIR THEORY OF INTEGRATION

As reference for this section Tao (2013) might serve.

In mathematics, a measure is a formal generalization of concepts such as length, area and volume.

More specifically, we are here interested in assigning a generalized notion of length (or mass) to subsets of the real half-line

$$\mathbb{R}_{\geq 0} = [0, \infty).$$

These sets will turn out to be elements of a so called σ -Algebra; i.e. a set Σ of sets for which

- $\emptyset, \mathbb{R}_{\geq 0} \in \Sigma$
- $A, B \in \Sigma \Rightarrow A \cap B \in \Sigma$
- $A, B \in \Sigma \Rightarrow A \setminus B \in \Sigma$
- $A, B \in \Sigma \Rightarrow A \cup B \in \Sigma$.

486 We now take $\Sigma_{\mathbb{R}_{\geq 0}}$ to be the smallest such set of sets Σ that contains all open intervals.

487 A complex measure then is a set-function that assigns to each set in $\Sigma_{\mathbb{R}_{\geq 0}}$ a complex number in a
488 certain way:

489 **Definition B.1.** A complex measure μ on $\mathbb{R}_{\geq 0}$ is a complex valued function $\mu : \Sigma_{\mathbb{R}_{\geq 0}} \rightarrow \mathbb{C}$ satisfying

$$491 \mu \left(\bigcup_n A_n \right) = \sum_n \mu(A_n)$$

492 for any countable (potentially infinite) collection of sets in $\Sigma_{\mathbb{R}_{\geq 0}}$ which are pairwise disjoint.

493 Let us provide some examples:

494 **Example B.2.** The prototypical example of a measure is the standard Lebesgue measure that assigns
495 to any interval (a, b) the length $\mu_{Leb}((a, b)) = |a - b|$ ($a, b \in \mathbb{R}_{\geq 0}$).

500 **Example B.3.** Alternatively, we might consider the Dirac measure $\mu_{\delta_{t_0}}$, which assigns the value
501 $\mu_{\delta_{t_0}}((a, b)) = 1$ to any interval (a, b) containing t_0 (i.e. $t_0 \in (a, b)$). Otherwise it assigns the value
502 $\mu_{\delta_{t_0}}((a, b)) = 0$ if $t_0 \notin (a, b)$.

503 **Example B.4.** Every integrable function $\hat{\psi} : \mathbb{R}_{\geq 0} \rightarrow \mathbb{C}$ defines a complex measure via $\mu_{\hat{\psi}}((a, b)) =$
504 $\int_a^b \hat{\psi}(t) dt$.

505 Hence we may think of **measures as generalizations of functions.**

506 Any given measure on $\mathbb{R}_{\geq 0}$ defines a unique way of integrating (known as Lebesgue inte-
507 gration) a function f defined on $\mathbb{R}_{\geq 0}$. This proceeds by approximating any function f via a weighted
508 sequence of indicator functions (with $A \in \Sigma_{\mathbb{R}_{\geq 0}}$ a set)

$$509 \chi_A(t) = \begin{cases} 1 & ; t \in A \\ 0 & ; t \notin A \end{cases}$$

510 as

$$511 f(t) \approx f_n(t) := \sum_k a_k^n \chi_{A_k}(t).$$

512 with $a_k \in \mathbb{C}$. For these functions, one then sets

$$513 \int_{\mathbb{R}_{\geq 0}} f_n d\mu \equiv \sum_k a_k^n \cdot \mu(A_k).$$

514 Since we have $\lim_{n \rightarrow \infty} f_n = f$, one then simply sets

$$515 \int_{\mathbb{R}_{\geq 0}} f d\mu \equiv \lim_{n \rightarrow \infty} \int_{\mathbb{R}_{\geq 0}} f_n d\mu.$$

516 **Example B.5.** For the prototypical example of the standard Lebesgue measure, this process simply
517 yields

$$518 \int_{\mathbb{R}_{\geq 0}} f(t) d\mu_{Leb}(t) = \int_0^{\infty} f(t) dt.$$

519 **Example B.6.** For the Dirac measure $\mu_{\delta_{t_0}}$, the above process yields

$$520 \int_{\mathbb{R}_{\geq 0}} f(t) d\mu_{\delta_{t_0}}(t) = f(t_0)$$

521 **Example B.7.** For measures arising from integrable functions $\hat{\psi} : \mathbb{R}_{\geq 0} \rightarrow \mathbb{C}$ as $\mu_{\hat{\psi}}((a, b)) =$
522 $\int_a^b \hat{\psi}(t) dt$, we find

$$523 \int_{\mathbb{R}_{\geq 0}} f(t) d\mu_{\hat{\psi}} = \int_0^{\infty} \hat{\psi}(t) f(t) dt.$$

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

B.2 LAPLACE TRANSFORMS

We say a complex valued measure μ is finite if we have

$$\int_{\mathbb{R}_{\geq 0}} d|\mu|(t) < \infty.$$

Here the measure $|\mu|$ arises from the original measure μ via

$$|\mu|((a, b)) \equiv |\mu((a, b))|.$$

For any such finite measure μ we may define its Laplace transform as

$$\psi_{\mu}(z) := \int_{\mathbb{R}_{\geq 0}} e^{-tz} d\mu(t).$$

This function f_{μ} is well defined for z in the right hemisphere

$$\mathbb{C}_R := \{z \in \mathbb{C} : \operatorname{Re}(z) \geq 0\}.$$

of the complex plane \mathbb{C} , since there we have

$$\begin{aligned} |\psi_{\mu}(z)| &= \left| \int_{\mathbb{R}_{\geq 0}} e^{-tz} d\mu(t) \right| \\ &\leq \int_{\mathbb{R}_{\geq 0}} |e^{-tz}| d|\mu|(t) \\ &\leq \int_{\mathbb{R}_{\geq 0}} d|\mu|(t) < \infty. \end{aligned}$$

Example B.8. For the Dirac measure $\mu_{\delta_{t_0}}$, we have

$$\psi_{\mu_{\delta_{t_0}}}(z) = e^{-t_0 z}.$$

Example B.9. For any integrable function $\hat{\psi}$, we have

$$\psi(z) \equiv \int_{\mathbb{R}_{\geq 0}} e^{-tz} d\mu_{\hat{\psi}} = \int_0^{\infty} \hat{\psi}(t) e^{-tz} dt.$$

More specifically, if the integrable function is given as $\hat{\psi}_k := (-t)^{k-1} e^{-\lambda t}$ (with $\operatorname{Re}(\lambda) > 0$), then $\psi_k(z) = (z + \lambda)^{-k}$:

Example B.10. If $\hat{\psi}_k := (-t)^{k-1} e^{-\lambda t}$ yields $\psi_k(z) = (z + \lambda)^{-k}$, then

$$\psi_k(z) = (z + \lambda)^{-k}.$$

For $k = 1$, this can be seen from

$$\int_0^{\infty} e^{-tz} e^{-\lambda t} dt = -\frac{1}{z + \lambda} e^{-(z+\lambda)t} \Big|_0^{\infty}.$$

For $k > 1$, the claim follows from differentiating the above expression with respect to z . Note that the functions $\psi_k(z) = (z + \lambda)^{-k}$ are also defined if $\operatorname{Re}(z) \leq 0$, as long as $z \neq -\lambda$.

Using the function ψ_k of the examples above, a wide class of functions may be parametrized

Theorem B.11. Let $f : \mathbb{R}_{\geq 0} \rightarrow 0$ be any function with $\lim_{x \rightarrow \infty} f(x) = 0$. Then for any $\epsilon > 0$, there is a function

$$h(x) = \sum_k \theta_k \psi_k(x)$$

for which

$$\sup_{x \in [0, \infty)} |f(x) - h(x)| < \epsilon.$$

Here the basis functions $\{\psi_k\}$ may either be chosen as $\psi_k(z) = (z + \lambda)^{-k}$ or $\psi_k(x) = e^{-(kt_0)x}$ for any $t_0 > 0$.

Proof. This is a direct consequence of the Weierstrass approximation theorem. \square

B.3 GLOBAL LAPLACIAN PROPAGATION MATRICES

A Global Laplacian Propagation matrix is then constructed by applying a function ψ arising as a Laplace transform to a graph Laplacian L . The resulting filter matrix $\psi(L) \in \mathbb{R}^{N \times N}$ acts on scalar graph signals $x \in \mathbb{R}^N$ via matrix multiplication; sending x to $\psi(L) \cdot x$:

$$x \mapsto \psi(L) \cdot x$$

C PROOFS RELATED TO GENERALIZATION ABILITY

C.1 GENERALIZATION ABILITY OF GLOBAL LAPLACIAN PROPAGATION MATRICES

In this section, we establish the generalization ability of global Laplacian propagation matrices.

Theorem C.1. *We have that $\|\psi(L) - J^\uparrow \psi(\underline{L}) J^\downarrow\| \leq \int_0^\infty |\hat{\psi}(t)| \eta(t) dt$ holds true.*

Proof. We start by proving the first claim. To this end, we note

$$\begin{aligned} \|\psi(L) - J^\downarrow \psi(\underline{L}) J^\uparrow\| &= \left\| \int_{\mathbb{R}_{\geq 0}} [e^{-tL} - J^\uparrow e^{-t\underline{L}} J^\downarrow] d\mu_{\hat{\psi}} \right\| \\ &\leq \int_{\mathbb{R}_{\geq 0}} \|e^{-tL} - J^\uparrow e^{-t\underline{L}} J^\downarrow\| d|\mu|_{\hat{\psi}} \end{aligned}$$

Using the notation for generalized fundtions, we have $d|\mu|_{\hat{\psi}}(t) = |\hat{\psi}(t)| dt$ and hence

$$\begin{aligned} \|\psi(L) - J^\downarrow \psi(\underline{L}) J^\uparrow\| &= \left\| \int_{\mathbb{R}_{\geq 0}} [e^{-tL} - J^\uparrow e^{-t\underline{L}} J^\downarrow] d\mu_{\hat{\psi}} \right\| \\ &\leq \int_{\mathbb{R}_{\geq 0}} \|e^{-tL} - J^\uparrow e^{-t\underline{L}} J^\downarrow\| |\hat{\psi}(t)| dt. \end{aligned}$$

□

Thus if $\eta(t) \equiv \|e^{-tL} - J^\uparrow e^{-t\underline{L}} J^\downarrow\| \approx 0$ on the support of $\hat{\psi}$, we also have $\|\psi(L) - J^\uparrow \psi(\underline{L}) J^\downarrow\| \approx 0$. In this case, propagation as implemented via $\psi(L)$ is essentially the same as propagation via $J^\downarrow \psi(\underline{L}) J^\uparrow$.

C.2 GENERALIZATION AND STABILITY WHEN $\|L - \tilde{L}\| \ll 1$

In this section we prove in addition to results in the main body of the paper also stability and generalization ability in the setting where for the Laplacians L, \tilde{L} of two graphs G, \tilde{G} defined on a common node set we have $\|L - \tilde{L}\| \ll 1$ (as opposed to the setting where one graph is a coarser version of another). We denote the collection of weight matrices by \mathcal{W} , the collection of biases by \mathcal{B} and the (collection of) utilized global Laplacian propagation matrices used in the update rule " $X \mapsto \sum_k \psi_k(L) X W_k$ " as Ψ . We denote the network by $\Phi_{\mathcal{W}, \mathcal{B}, \Psi}$ and write the generated embeddings for the node feature matrix X as $\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(X)$. With this, we have:

Theorem C.2. *Let $\Phi_{\mathcal{W}, \mathcal{B}, \Psi}$ be a K -layer deep graph convolutional architecture. Assume in each layer $1 \leq \ell \leq K$ that $\sum_i \|W_i^\ell\| \leq W$ and $\|B^\ell\| \leq B$. Choose $C \geq \|\Psi_i(L)\|$ ($\forall i \in I$) and w.l.o.g. assume $CW > 1$. With this, we have with $\delta = \max_{i \in I} \{\|\Psi_i(L) - \Psi_i(\tilde{L})\|\}$ that*

$$\|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\tilde{L}, X)\| \leq \left[K \cdot C^K W^{K-1} \cdot \left(\|X\| + \frac{1}{CW - 1} B \right) \right] \cdot \delta.$$

648 *Proof.* For simplicity in notation, let us denote the hidden representations in the network correspond-
 649 ing to \tilde{L} by X^ℓ . With this, we note:

$$\begin{aligned}
 650 \quad \|X^K - \tilde{X}^K\| &\leq \sum_{i \in I} \|\psi_i(L) - \psi_i(\tilde{L})\| \cdot \|X^{K-1}\| \cdot \|W_i^K\| + \sum_{i \in I} \|\psi_i(\tilde{L})\| \cdot \|\tilde{X}^{K-1} - X^{K-1}\| \cdot \|W_i^K\| \\
 651 \quad &\leq \delta W \|X^{K-1}\| + CW \|\tilde{X}^{K-1} - X^{K-1}\| \\
 652 \quad &\leq \delta W \|X^{K-1}\| + CW \delta \|X^{K-2}\| + (CW)^2 \|\tilde{X}^{K-1} - X^{K-1}\| \\
 653 \quad &\leq \frac{\delta}{C} \cdot \left(\sum_{\ell=1}^K (CW)^\ell \|X^{K-\ell}\| \right) \\
 654 \quad &= \frac{\delta}{C} \cdot \left(\sum_{j=0}^{K-1} (CW)^{K-j} \|X^j\| \right) \\
 655 \quad &
 \end{aligned}$$

656 Hence we need to bound the quantity $\|X^j\|$ in terms of C, W, B and X .

657 We have

$$\begin{aligned}
 658 \quad \|X^j\| &\leq \sum_i \|\psi_i(L)\| \cdot \|X^{j-1}\| \cdot \|W_i^j\| + \|B^j\| \\
 659 \quad &\leq CW \|X^{j-1}\| + B \\
 660 \quad &\leq (CW)^2 \|X^{j-2}\| + CW B + B \\
 661 \quad &\leq B \left(\sum_{k=0}^{j-1} (CW)^k \right) + (CW)^j \|X\| \\
 662 \quad &= \begin{cases} B \frac{(CW)^j - 1}{CW - 1} + (CW)^j \|X\| & ; CW \neq 1 \\ jB + \|X\| & ; CW = 1 \end{cases} . \\
 663 \quad &
 \end{aligned}$$

664 For the case $CW = 1$, we thus find

$$\begin{aligned}
 665 \quad \|X^K - \tilde{X}^K\| &\leq \frac{\delta}{C} \cdot \left(\sum_{j=0}^{K-1} (jB + \|X\|) \right) \\
 666 \quad &= \frac{\delta}{C} \cdot \left(K\|X\| + B \frac{K(K-1)}{2} \right) . \\
 667 \quad &
 \end{aligned}$$

668 For the case $CW \neq 1$, we find

$$\|X^K - \tilde{X}^K\| \leq \frac{\delta}{C} \cdot \left(\sum_{j=0}^{K-1} (CW)^{K-j} \left[B \frac{(CW)^j - 1}{CW - 1} + (CW)^j \|X\| \right] \right)$$

669 For $CW > 1$, we may further estimate this as

$$\begin{aligned}
 670 \quad \|X^K - \tilde{X}^K\| &\leq \frac{\delta}{C} \cdot \left(\sum_{j=0}^{K-1} (CW)^{K-j} \left[B \frac{(CW)^j - 1}{CW - 1} + (CW)^j \|X\| \right] \right) \\
 671 \quad &\leq \delta \cdot \frac{K(CW)^K}{C} \left[\frac{B}{CW - 1} + \|X\| \right] . \\
 672 \quad &
 \end{aligned}$$

673 This proves the claim. \square

674 C.3 PROOF OF CONVERGENCE IN THE SENSE OF SECTION 3

675 The result in Section 3 is concerned with the graph-level setting; i.e. the setting where entire graphs
 676 are embedded into latent spaces. Before proving this result, we first prove a corresponding result for
 677 the node-level, where individual nodes in a graph are embedded. We will then use this node-level

702 result (Theorem C.3 below) to prove the graph-level Convergence result in section 3.
703

704 In the node-level setting, we start by considering initial node-features X on G . We then fix
705 a graph neural network Φ based on global Laplacian propagation schemes and consider two ways of
706 generating embeddings on the graph G : On the one hand, we may simply generate embeddings with
707 the network Φ on G . On the other hand, we may also project the node feature matrix X to \underline{G} via
708 J^\downarrow , apply the network Φ to the matrix $J^\downarrow X$ on \underline{G} and then finally interpolate the generated node
709 embeddings back to G via J^\uparrow .

710 The following result bounds the difference between these two respective node embeddings generated
711 on the same graph.

712 **Theorem C.3.** *Let $\Phi_{\mathcal{W}, \mathcal{B}, \Psi}$ be a K -layer deep Global-Laplacian-Propagation-based network. Assume $\sum_{i \in I} \|W_i^\ell\| \leq W$ and bound bias matrices in layer ℓ as $\|B^\ell\| \leq B$. Choose $C \geq \|\Psi_i(L)\|$ ($i \in I$) and w.l.o.g. assume $CW > 1$ (which can always be satisfied by choosing C large enough). Assume $\rho(J^\uparrow X) = J^\uparrow \rho(X)$ and if biases are enabled, assume $J^\uparrow \mathbf{1}_{\underline{G}} = \mathbf{1}_G$. Set $\max_{i \in I} \{\|\psi_i(L) - J^\uparrow \psi_i(\underline{L}) J^\downarrow\|\} = \delta_1$ and define $\delta_2 = \max_{i \in I} \{\|\psi_i(L^\uparrow)[J^\downarrow J^\uparrow - Id_{\underline{G}}]\|\}$. With this, we have that*

$$713 \|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X)\| \leq \left[K \cdot C^K W^{K-1} \cdot \left(\|X\| + \frac{1}{CW-1} B \right) \right] \cdot (\delta_1 + \delta_2). \\ 720$$

721 It should be noted that the result above is more general than the setting considered in Section 3. In
722 the setting considered in Section 3 we have $J^\downarrow J^\uparrow = Id_{\underline{G}}$ (in addition to $\rho(J^\uparrow X) = J^\uparrow \rho(X)$). There
723 we thus automatically have $\delta_2 = 0$.
724

725 *Proof.* Let us define

$$726 \underline{X} := J^\downarrow X.$$

727 Let us further use the notation $\underline{\psi}_i := \psi_i(\underline{L})$ and $\psi_i := \psi_i(L)$.
728

729 Denote by X^ℓ and \underline{X}^ℓ the (hidden) feature matrices generated in layer ℓ for networks based on ψ_i
730 and $\underline{\psi}_i$ respectively: I.e. we have
731

$$732 X^\ell = \rho \left(\sum_{i \in I} \psi_i X^{\ell-1} W_i^\ell + B^\ell \right)$$

733 and
734

$$735 \underline{X}^\ell = \rho \left(\sum_{i \in I} \underline{\psi}_i \underline{X}^{\ell-1} W_i^\ell + \underline{B}^\ell \right). \\ 736$$

737 We then have
738

$$739 \begin{aligned} & \|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X)\| \\ &= \|X^K - J^\uparrow \underline{X}^K\| \\ &= \left\| \rho \left(\sum_{i \in I} \psi_i X^{K-1} W_i^K + B^K \right) - J^\uparrow \rho \left(\sum_{i \in I} \underline{\psi}_i \underline{X}^{K-1} W_i^K + \underline{B}^L \right) \right\| \\ &= \left\| \rho \left(\sum_{i \in I} \psi_i X^{K-1} W_i^K + B^K \right) - \rho \left(J \sum_{i \in I} \underline{\psi}_i \underline{X}^{K-1} W_i^K + B^L \right) \right\| \end{aligned} \\ 740 \\ 741 \\ 742 \\ 743$$

744 Here we used the assumption that ρ and J commute. In fact since $\text{ReLU}(\cdot)$ maps positive entries
745 to positive entries and acts pointwise, it commutes with J^\uparrow . We also made use of the assumption
755 $J^\uparrow \mathbf{1}_{\underline{G}} = \mathbf{1}_G$ when dealing with biases .

Using the fact that $\rho(\cdot)$ is 1-Lipschitz-continuous, we can establish

$$\begin{aligned} & \|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, JX)\| \\ & \leq \left\| \rho \left(\sum_{i \in I} \psi_i X^{K-1} W_i^K + B^K \right) - \rho \left(J^\uparrow \sum_{i \in I} \underline{\psi}_i \underline{X}^{K-1} W_i^K + B^K \right) \right\| \\ & \leq \left\| \sum_{i \in I} \psi_i X^{K-1} W_i^K + B^K - J^\uparrow \sum_{i \in I} \underline{\psi}_i \underline{X}^{K-1} W_i^K + B^K \right\|. \end{aligned}$$

Using the assumption that $\|\psi[J^\downarrow J^\uparrow - Id_G]\| \leq \delta_2$, we have

$$\begin{aligned} & \|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, JX)\| \\ & \leq \left\| \sum_{i \in I} \psi_i X^{K-1} W_i^K - \sum_{i \in I} (J^\uparrow \underline{\psi}_i J) J^\uparrow \underline{X}^{K-1} W_i^K \right\| + \left\| \sum_{i \in I} J^\uparrow \underline{\psi}_i [Id_G - J^\downarrow J^\uparrow] \underline{X}^{K-1} W_i^K \right\| \\ & \leq \left\| \sum_{i \in I} \psi_i X^{K-1} W_i^K - \sum_{i \in I} (J^\uparrow \underline{\psi}_i J) J^\uparrow \underline{X}^{K-1} W_i^K \right\| + \delta_2 \cdot \left\| \sum_{i \in I} \underline{X}^{K-1} W_i^K \right\| \\ & \leq \left\| \sum_{i \in I} \psi_i X^{K-1} W_i^K - \sum_{i \in I} (J^\uparrow \underline{\psi}_i J^\downarrow) J^\uparrow \underline{X}^{K-1} W_i^K \right\| + \delta_2 \cdot \|\underline{X}^{K-1}\| \cdot W \end{aligned}$$

From this, we find (assuming $\|J^\uparrow\|, \|J^\downarrow\| \leq 1$ for notational simplicity (and which is true in the setting of Section 3)), that

$$\begin{aligned} & \|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, JX)\| \\ & \leq \left\| \sum_{i \in I} \psi_i X^{K-1} W_i^K - \sum_{i \in I} (J^\uparrow \underline{\psi}_i J^\downarrow) J^\uparrow \underline{X}^{K-1} W_i^K \right\| + \delta_2 \cdot \|\underline{X}^{K-1}\| \cdot W \\ & \leq \left\| \sum_{i \in I} (\psi_i - J^\uparrow \underline{\psi}_i J) X^{K-1} W_i^K \right\| + \sum_{i \in I} \|J^\uparrow \underline{\psi}_i J\| \cdot \|J^\uparrow \underline{X}^{K-1} - X^{K-1}\| \cdot \|W_i^K\| + \delta_2 \cdot \|\underline{X}^{K-1}\| \cdot W \\ & \leq \left\| \sum_{i \in I} (\psi_i - J^\uparrow \underline{\psi}_i J) X^{K-1} W_i^K \right\| + CW \cdot \|J^\uparrow \underline{X}^{K-1} - X^{K-1}\| + \delta_2 \cdot \|\underline{X}^{K-1}\| \cdot W \\ & \leq \sum_{i \in I} \left\| (\psi_i - J^\uparrow \underline{\psi}_i J) \right\| \cdot \|X^{K-1}\| \cdot \|W_i^K\| + CW \cdot \|J^\uparrow \underline{X}^{K-1} - X^{K-1}\| + \delta_2 \cdot \|\underline{X}^{K-1}\| \cdot W \\ & \leq \delta_1 \cdot \|X^{K-1}\| W + CW \cdot \|J^\uparrow \underline{X}^{K-1} - X^{K-1}\| + \delta_2 \cdot \|\underline{X}^{K-1}\| \cdot W \end{aligned}$$

Arguing as in the proof of Appendix C.2 then yields the claim. \square

Let us move from the node-level to the graph-level. We first specify how graph-level latent embeddings arise:

Definition C.4. We aggregate embeddings $X \in \mathbb{R}^{N \times F}$ of individual nodes to graph-embeddings $\Omega(X) \in \mathbb{R}^F$ as $\Omega(X)_j = \sum_{i=1}^N |X_{ij}| \cdot \mu_i$. Here $\{\mu_i\}_i$ is the set of node-weights.

In a social network, a node weight $\mu_i = 1$ might e.g. signify that node i represents a single user. A weight $\mu_j > 1$ would indicate that node j represents a group of users.

Given such an aggregation of node embeddings into latent-embeddings of entire graphs, we may then relegate graph-level transferability back to node-level transferability:

Theorem C.5. Assuming $\Omega(\underline{X}) = \Omega(J^\uparrow X)$, we have in the setting of Theorem C.3 that $\|\Omega \circ \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - \Omega \circ \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, J^\downarrow X)\| \leq \|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, J^\downarrow X)\|$.

810 *Proof.* We note

$$\begin{aligned}
811 & \|\Omega \circ \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X) - \Omega \circ \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X)\| \\
812 & = \|\Omega(\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X)) - \Omega(\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X))\| \\
813 & = \|\Omega(\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L, X)) - \Omega(J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X))\|.
\end{aligned}$$

816 To prove the claim from here, we only have to note that the aggregation method Ω as defined in
817 Definition C.5 above is 1-Lipschitz (as a consequence of the reverse triangle inequality). The proof
818 for the bidirectional setting proceeds analogously. \square

820 This result then proves the continuity result of Section 3. Indeed: In the notation of Section 3, we
821 have $F_\omega = \Omega(\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L_\omega, X))$ and $\underline{F} = \Omega(\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X))$. Thus we have

$$822 \quad \|F_\omega - \underline{F}\| = \|\Omega \circ \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L_\omega, X) - \Omega \circ \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X)\| \leq \|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L_\omega, X) - J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X)\|.$$

824 By Theorem C.3 and the fact that $[Id_{\underline{G}} - J^\uparrow J^\downarrow] = 0$, we have

$$825 \quad \|\Phi_{\mathcal{W}, \mathcal{B}, \Psi}(L_\omega, X) - J^\uparrow \Phi_{\mathcal{W}, \mathcal{B}, \Psi}(\underline{L}, J^\downarrow X)\| \lesssim \max_k \{\|\psi_k(L_\omega) - J^\uparrow \psi_k(\underline{L}) J^\downarrow\|\},$$

828 with " \lesssim " as per usual "denoting smaller than, up to a positive multiplicative constant".

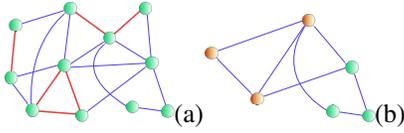
829 Finally Theorem C.1 implies

$$831 \quad \|\psi_k(L_\omega) - J^\uparrow \psi_k(\underline{L}) J^\downarrow\| \leq \int_0^\infty |\hat{\psi}_k(t)| \eta(t) dt = \int_{\mathbb{R}_{\geq 0}} \|e^{-tL_\omega} - J^\uparrow e^{-t\underline{L}} J^\downarrow\| |\hat{\psi}_k(t)| dt.$$

834 Thus upon combining these steps and noting that $\eta_\omega \rightarrow 0$ by (1), the convergence result of Section 3
835 is indeed proved.

837 D ADDITIONAL EXPERIMENTAL CONSIDERATIONS

840 Collapsing strongly connected clusters: Intuition and exact Definitions



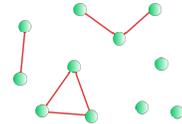
848 Figure 4: (a) G (strongly connected)
849 clusters in red (b) Coarse grained \underline{G}

850 in G_{cluster} (c.f. Fig 5). Edges \mathcal{E} are given by elements $(R, P) \in \underline{G} \times \underline{G}$ with non-zero accumu-
851 lated edge weight $\underline{W}_{RP} = \sum_{r \in R} \sum_{p \in P} W_{rp}$. Node weights in \underline{G} are defined accordingly by
852 aggregating as $\underline{\mu}_R = \sum_{r \in R} \mu_r$. To compare signals on these two graphs, we
853 define intertwining operators J^\downarrow, J^\uparrow transferring information between G and \underline{G} :
854 Let x be a scalar graph signal and let $\mathbf{1}_R$ be the vector that has 1 as entry for
855 nodes $r \in R$ and is zero otherwise. Denote by u_R the entry of u at node $R \in \underline{G}$.
856 Projection J^\downarrow is then defined component-wise by evaluation at node $R \in \underline{G}$ as the
857 average of x over R : $(J^\downarrow x)_R = \langle \mathbf{1}_R, x \rangle / \underline{\mu}_R$. Going in the opposite direction,
858 interpolation is defined as $J^\uparrow u = \sum_{R \in \underline{G}} u_R \cdot \mathbf{1}_R$.

860 In this setting, we have (c.f. Appendix A) that

$$861 \quad \|e^{-tL} - J^\uparrow e^{-t\underline{L}} J^\downarrow\| \lesssim 1/w_{\text{high}}^{\min} \text{ for any } t > 0.$$

862 Here $w_{\text{high}}^{\min} \gg 1$ denotes the minimal edge weight inside the strongly connected clusters in G .



863 Figure 5: G_{cluster}

Dataset: The dataset we consider is the **QM7** dataset, introduced in Blum & Raymond (2009); Rupp et al. (2012). This dataset contains descriptions of 7165 organic molecules, each with up to seven heavy atoms, with all non-hydrogen atoms being considered heavy. A molecule is represented by its Coulomb matrix C^{Cmb} , whose off-diagonal elements

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

correspond to the Coulomb-repulsion between atoms i and j . We discard diagonal entries of Coulomb matrices; which would encode a polynomial fit of atomic energies to nuclear charge Rupp et al. (2012).

For each atom in any given molecular graph, the individual Cartesian coordinates R_i and the atomic charge Z_i are (in principle) also accessible individually. To each molecule an atomization energy - calculated via density functional theory - is associated. The objective is to predict this quantity. The performance metric is mean absolute error. Numerically, atomization energies are negative numbers in the range -600 to -2200 . The associated unit is $[kcal/mol]$.

Details on collapsing procedure as applied to QM7: Again, we make use of the QM7 dataset Rupp et al. (2012) and its Coulomb matrix description

$$C_{ij}^{\text{Cmb}} = \frac{Z_i Z_j}{|R_i - R_j|} \quad (4)$$

of molecules. We modify (all) molecular graphs in QM7 by deflecting hydrogen atoms (H) out of their equilibrium positions towards the respective nearest heavy atom. This is possible since the QM7 dataset also contains the Cartesian coordinates of individual atoms. Edge weights between heavy atoms then remain the same, while Coulomb repulsions between H-atoms and respective nearest heavy atom increasingly diverge; as is evident from (4).

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

Mathematically, \underline{G} is obtained by removing all nodes corresponding to H-atoms from G , while adding the corresponding charges $Z_H = 1$ to the node-weights of the respective nearest heavy atom. Charges in (4) are modified similarly to generate the weight matrix \underline{W} .

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

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

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

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

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

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

The feature matrix

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

is a single row-vector given as

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

918 **Experimental Setup:** We randomly select 1500 molecules for testing and train on the remaining
 919 graphs. On QM7 we run experiments for 23 different random seeds and report mean and
 920 standard deviation. All experiments were performed on a single NVIDIA Quadro RTX 8000 graphics
 921 card.

922
 923 **Additional details on training and models:** Typical GNN models are divided into **standard**
 924 architectures (GCN (Kipf & Welling, 2017), ChebNet (Defferrard et al., 2016), ARMA (Bianchi
 925 et al., 2019), BernNet (He et al., 2021), GATv2 (Brody et al., 2022)) and **multi-scale** architectures
 926 (PushNet (Busch et al., 2020), UFGNet (Zheng et al., 2021), Lanczos (Liao et al., 2019)). Apart
 927 from UFGNet (already acting as a **pooling** layer) we also consider self-attention-pooling (Lee et al.,
 928 2019); both acting on the final layer (SAG) and as acting on the output of each individual layer, with
 929 resulting layer-wise features concatenated to produce the final embedding (SAG-M). All considered
 930 convolutional layers are incorporated into a two layer deep and fully connected graph convolutional
 931 architecture. In each hidden layer, we set the width (i.e. the hidden feature dimension) to

$$F_1 = F_2 = 64.$$

932
 933 For BernNet, we set the polynomial order to $K = 3$ to combat appearing numerical instabilities.
 934 ARMA is set to $K = 2$ and $T = 1$. ChebNet uses $K = 2$. Lanczos uses 20 Lanczos iterations, as
 935 proposed in the original paper (Liao et al., 2019). UFGNet uses Haar wavelets. For all baselines,
 936 the standard mean-aggregation scheme is employed after the graph-convolutional layers to generate
 937 graph level features. Finally, predictions are generated via an MLP.

938 For the **resolvent** based global Laplacian propagation architecture, we set $\lambda = 1$ and build filters
 939 using the $k = 1$ and $= 2$ matrices in $\Psi^{\text{Res}} = \{(z + \lambda)^{-k}\}_{k \in \mathbb{N}}$.

940
 941 For the **based global Laplacian propagation architecture**, based global Laplacian propagation
 942 architecture, we set $t_0 = 1$ and build filters using the $k = 1$ and $= 2$ matrices in $\Psi^{\text{Exp}} =$
 943 $\{e^{-(kt_0)z}\}_{k \in \mathbb{N}}$.

944 As aggregation, we employ the graph level feature aggregation scheme introduced in Definition C.4
 945 with node weights set to atomic charges of individual atoms. Predictions are then generated via a
 946 final MLP with the same specifications as the one used for baselines.

947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971