

---

# Graph Neural Networks as Gradient Flows

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Dynamical systems minimizing an energy are ubiquitous in geometry and physics. We propose a gradient flow framework for GNNs where the equations follow the direction of steepest descent of a learnable energy. This approach allows to analyse the GNN evolution from a multi-particle perspective as learning attractive and repulsive forces in feature space via the positive and negative eigenvalues of a symmetric ‘channel-mixing’ matrix. We perform spectral analysis of the solutions and conclude that gradient flow graph convolutional models can induce a dynamics dominated by the graph high frequencies, which is desirable for heterophilic datasets. We also describe structural constraints on common GNN architectures allowing to interpret them as gradient flows. We perform thorough ablation studies corroborating our theoretical analysis and show competitive performance of simple and lightweight models on real-world homophilic and heterophilic datasets.

## 1 Introduction and motivations

Graph neural networks (GNNs) [38, 20, 21, 36, 7, 15, 27] and in particular their Message Passing formulation (MPNN) [19] have become the standard ML tool for dealing with different types of relations and interactions, ranging from social networks to particle physics and drug design. One of the often cited drawbacks of traditional GNN models is their poor ‘explainability’, making it hard to know why and how they make certain predictions [46, 47], and in which situations they may work and when they would fail. Limitations of GNNs that have attracted attention are over-smoothing [29, 30, 8], over-squashing and bottlenecks [1, 40], and performance on heterophilic data [31, 51, 13, 4, 45] – where adjacent nodes usually have different labels.

**Contributions.** We propose a *Gradient Flow Framework* (GRAFF) where the GNN equations follow the direction of steepest descent of a *learnable energy*. Thanks to this framework we can (i) interpret GNNs as a multi-particle dynamics where the learned parameters determine pairwise attractive and repulsive potentials in the feature space. This sheds light on how GNNs can adapt to heterophily and explains their performance and the smoothness of the prediction. (ii) GRAFF leads to residual convolutional models where the *channel-mixing*  $\mathbf{W}$  is performed by a shared symmetric bilinear form inducing attraction and repulsion via its positive and negative eigenvalues, respectively. We theoretically investigate the interaction of the graph spectrum with the spectrum of the channel-mixing, proving that if there is more mass on the negative eigenvalues of  $\mathbf{W}$ , then the dynamics is dominated by the graph-high frequencies, which could be desirable on heterophilic graphs. We also extend results of [29, 30, 8] by showing that when we drop the residual connection intrinsic to the gradient flow framework,

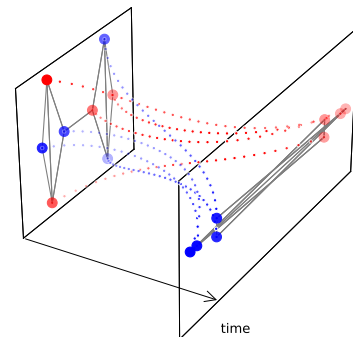


Figure 1: GRAFF dynamics: attractive and repulsive forces lead to a non-smoothing process able to separate labels.

graph convolutional models always induce a low-frequency dominated dynamics *independent* of the sign and magnitude of the spectrum of the channel-mixing. We also discuss how simple choices make common architectures fit GRAFF and conduct thorough ablation studies to corroborate the theoretical analysis on the role of the spectrum of  $\mathbf{W}$ . (iii) We crystallize *an instance* of our framework into a linear, residual, convolutional model that achieves competitive performance on homophilic and heterophilic real world graphs whilst being faster than GCN.

**Related work.** Our analysis is related to studying GNNs as filters on the graph spectrum [15, 24, 2, 25] and over-smoothing [29, 30, 8, 50] and partly adopts techniques similar to [30]. The key difference is that we also consider the spectrum of the ‘channel-mixing’ matrix. The concept of gradient flows has been a standard tool in physics and geometry [16], from which they were adopted for image processing [26], and recently used in ML [35] for the analysis of Transformers [41] – see also [18] for discussion of loss landscapes. Our continuous-time evolution equations follows the spirit of Neural ODES [22, 12, 3] and the study of GNNs as continuous dynamical systems [44, 10, 17, 9].

**Outline.** In Section 2, we review the continuous and discrete Dirichlet energy and the associated gradient flow framework. We formalize the notion of over-smoothing and low(high)-frequency-dominated dynamics to investigate GNNs and study the dominant components in their evolution. We extend the graph Dirichlet energy to allow for a non-trivial norm for the feature edge-gradient. This leads to gradient flow equations that diffuse the features and over-smooth in the limit. Accordingly, in Section 3 we introduce a more general energy with a symmetric channel-mixing matrix  $\mathbf{W}$  giving rise to attractive and repulsive pairwise terms via its positive and negative eigenvalues and show that the negative spectrum can induce high-frequency-dominant dynamics. In Section 4 we first compare with continuous GNN models and then discretize the equations and provide a ‘recipe’ for making standard GNN architectures fit a gradient flow framework. We adapt the spectral analysis to discrete-time showing that gradient flow convolutional models *can* generate a dynamics dominated by the high frequencies via the negative eigenvalues of  $\mathbf{W}$  while this is impossible if we drop the residual connection. In Section 5 we corroborate our theoretical analysis on the role of the spectrum of  $\mathbf{W}$  via ablation studies on graphs with varying homophily. Experiments on real world datasets show a competitive performance of our model despite its simplicity and reduced number of parameters.

## 2 Gradient-flow formalism

**Notations adopted throughout the paper.** Let  $G = (V, E)$  be an *undirected* graph with  $n$  nodes. We denote by  $\mathbf{F} \in \mathbb{R}^{n \times d}$  the matrix of  $d$ -dimensional node features, by  $\mathbf{f}_i \in \mathbb{R}^d$  its  $i$ -th row (transposed), by  $\mathbf{f}^r \in \mathbb{R}^n$  its  $r$ -th column, and by  $\text{vec}(\mathbf{F}) \in \mathbb{R}^{nd}$  the vectorization of  $\mathbf{F}$  obtained by stacking its columns. Given a symmetric matrix  $\mathbf{B}$ , we let  $\lambda_+^{\mathbf{B}}, \lambda_-^{\mathbf{B}}$  denote its most positive and negative eigenvalues, respectively, and  $\rho_{\mathbf{B}}$  be its *spectral radius*. If  $\mathbf{B} \succeq 0$ , then  $\text{gap}(\mathbf{B})$  denotes the *positive smallest eigenvalue* of  $\mathbf{B}$ .  $\dot{f}(t)$  denotes the temporal derivative,  $\otimes$  is the Kronecker product and ‘a.e.’ means *almost every* w.r.t. Lebesgue measure and usually refers to data in the complement of some lower dimensional subspace in  $\mathbb{R}^{n \times d}$ . Proofs and additional results appear in the Appendix.

**Starting point: a geometric parallelism.** To motivate a gradient-flow approach for GNNs, we start from the continuous case (see Appendix A.1 for details). Consider a smooth map  $f : \mathbb{R}^n \rightarrow (\mathbb{R}^d, h)$  with  $h$  a constant metric represented by  $\mathbf{H} \succeq 0$ . The *Dirichlet energy* of  $f$  is defined by

$$\mathcal{E}(f, h) = \frac{1}{2} \int_{\mathbb{R}^n} \|\nabla f\|_h^2 dx = \frac{1}{2} \sum_{q,r=1}^d \sum_{j=1}^n \int_{\mathbb{R}^n} h_{qr} \partial_j f^q \partial_j f^r(x) dx \quad (1)$$

and measures the ‘smoothness’ of  $f$ . A natural approach to find minimizers of  $\mathcal{E}$  - called *harmonic maps* - was introduced in [16] and consists in studying the **gradient flow** of  $\mathcal{E}$ , wherein a given map  $f(0) = f_0$  is evolved according to  $\dot{f}(t) = -\nabla_f \mathcal{E}(f(t))$ . These type of evolution equations have historically been the core of *variational* and *PDE-based image processing*; in particular, gradient flows of the Dirichlet energy were shown [26] to recover the Perona-Malik nonlinear diffusion [32].

**Motivation: GNNs for node-classification.** We wish to extend the gradient flow formalism to node classification on graphs. Assume we have a graph  $G$ , node-features  $\mathbf{F}_0$  and labels  $\{y_i\}$  on  $V_{\text{train}} \subset V$ , and that we want to predict the labels on  $V_{\text{test}} \subset V$ . A GNN typically evolves the features via some

parametric rule,  $\text{GNN}_\theta(\mathbf{G}, \mathbf{F}_0)$ , and uses a decoding map for the prediction  $y = \psi_{\text{DE}}(\text{GNN}_\theta(\mathbf{G}, \mathbf{F}_0))$ . In graph convolutional models [15, 27],  $\text{GNN}_\theta$  consists of two operations: applying a shared linear transformation to the features (**‘channel mixing’**) and propagating them along the edges of the graph (**‘diffusion’**). Our **goal** consists in studying when  $\text{GNN}_\theta$  is the *gradient flow* of some parametric class of energies  $\mathcal{E}_\theta : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}$ , which generalize the Dirichlet energy. This means that the parameters can be interpreted as ‘finding the right notion of smoothness’ for our task. We evolve the features by  $\dot{\mathbf{F}}(t) = -\nabla_{\mathbf{F}} \mathcal{E}_\theta(\mathbf{F}(t))$  with prediction  $y = \psi_{\text{DE}}(\mathbf{F}(T))$  for some optimal time  $T$ .

**Why a gradient flow?** Since  $\dot{\mathcal{E}}_\theta(\mathbf{F}(t)) = -\|\nabla_{\mathbf{F}} \mathcal{E}_\theta(\mathbf{F}(t))\|^2$ , the energy dissipates along the gradient flow. Accordingly, this framework allows to *explain the GNN dynamics* as flowing the node features in the direction of steepest descent of  $\mathcal{E}_\theta$ . Indeed, we find that parametrizing an energy leads to equations governed by attractive and repulsive forces that can be controlled via the spectrum of symmetric ‘channel-mixing’ matrices. This shows that by learning to distribute more mass over the negative (positive) eigenvalues of the channel-mixing, gradient flow models can generate dynamics dominated by the higher (respectively, lower) graph frequencies and hence tackle different homophily scenarios. The gradient flow framework also leads to sharing of the weights across layers (since we parametrize the *energy* rather than the *evolution equations*, as usually done in GNNs), allowing us to reduce the number of parameters without compromising performance (see Table 1).

**Analysis on graphs: preliminaries.** Given a *connected* graph  $\mathbf{G}$  with self-loops, its adjacency matrix  $\mathbf{A}$  is defined as  $a_{ij} = 1$  if  $(i, j) \in \mathbf{E}$  and zero otherwise. We let  $\mathbf{D} = \text{diag}(d_i)$  be the degree matrix and write  $\bar{\mathbf{A}} := \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ . Let  $\mathbf{F} \in \mathbb{R}^{n \times d}$  be the matrix representation of a signal. Its *graph gradient* is  $(\nabla \mathbf{F})_{ij} := \mathbf{f}_j / \sqrt{d_j} - \mathbf{f}_i / \sqrt{d_i}$ . We define the *Laplacian* as  $\Delta := -\frac{1}{2} \text{div } \nabla$  (the *divergence*  $\text{div}$  is the adjoint of  $\nabla$ ), represented by  $\Delta = \mathbf{I} - \bar{\mathbf{A}} \succeq 0$ . We refer to the eigenvalues of  $\Delta$  as *frequencies*: the lowest frequency is always 0 while the highest frequency is  $\rho_\Delta \leq 2$  [14]. As for the continuum case, the gradient allows to define a (*graph*) *Dirichlet energy* as [49]

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}) := \frac{1}{4} \sum_i \sum_{j:(i,j) \in \mathbf{E}} \|(\nabla \mathbf{F})_{ij}\|^2 \equiv \frac{1}{4} \sum_{(i,j) \in \mathbf{E}} \left\| \frac{\mathbf{f}_i}{\sqrt{d_i}} - \frac{\mathbf{f}_j}{\sqrt{d_j}} \right\|^2 = \frac{1}{2} \text{trace}(\mathbf{F}^\top \Delta \mathbf{F}), \quad (2)$$

where the extra  $\frac{1}{2}$  is for convenience. As for manifolds,  $\mathcal{E}^{\text{Dir}}$  measures smoothness. If we stack the columns of  $\mathbf{F}$  into  $\text{vec}(\mathbf{F}) \in \mathbb{R}^{nd}$ , the gradient flow of  $\mathcal{E}^{\text{Dir}}$  yields the *heat equation* on each channel:

$$\text{vec}(\dot{\mathbf{F}}(t)) = -\nabla_{\text{vec}(\mathbf{F})} \mathcal{E}^{\text{Dir}}(\text{vec}(\mathbf{F}(t))) = -(\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F}(t)) \iff \dot{\mathbf{f}}^r(t) = -\Delta \mathbf{f}^r(t), \quad (3)$$

for  $1 \leq r \leq d$ . Similarly to [8], we rely on  $\mathcal{E}^{\text{Dir}}$  to assess whether a given dynamics  $t \mapsto \mathbf{F}(t)$  is a smoothing process. A different choice of Laplacian  $\mathbf{L} = \mathbf{D} - \mathbf{A}$  with non-normalized adjacency induces the analogous Dirichlet energy  $\mathcal{E}_{\mathbf{L}}^{\text{Dir}}(\mathbf{F}) = \frac{1}{2} \text{trace}(\mathbf{F}^\top \mathbf{L} \mathbf{F})$ . Throughout this paper, we rely on the following definitions (see Appendix A.3 for further equivalent formulations and justifications):

**Definition 2.1.**  $\dot{\mathbf{F}}(t) = \text{GNN}_\theta(\mathbf{F}(t), t)$  initialized at  $\mathbf{F}(0)$  is *smoothing* if  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \leq C + \varphi(t)$ , with  $C$  a constant only depending on  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(0))$  and  $\dot{\varphi}(t) \leq 0$ . *Over-smoothing* occurs if either  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow 0$  or  $\mathcal{E}_{\mathbf{L}}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow 0$  for  $t \rightarrow \infty$ .

Our notion of ‘over-smoothing’ is a relaxed version of the definition in [34] – although in the linear case one always finds an *exponential decay* of  $\mathcal{E}^{\text{Dir}}$ . We note that  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow 0$  iff  $\Delta \mathbf{f}^r(t) \rightarrow \mathbf{0}$  for each column  $\mathbf{f}^r$ . As in [30], this corresponds to a loss of separation power along the solution where nodes with *equal degree* become indistinguishable since we converge to  $\ker(\Delta)$  (if we replaced  $\Delta$  with  $\mathbf{L}$  then we would not even be able to separate nodes with different degrees in the limit).

To motivate the next definition, consider  $\dot{\mathbf{F}}(t) = \bar{\mathbf{A}} \mathbf{F}(t)$ . Despite  $\|\mathbf{F}(t)\|$  being unbounded for a.e.  $\mathbf{F}(0)$ , the low-frequency components are growing the fastest and indeed  $\mathbf{F}(t)/\|\mathbf{F}(t)\| \rightarrow \mathbf{F}_\infty$  s.t.  $\Delta \mathbf{f}_\infty^r = \mathbf{0}$  for  $1 \leq r \leq d$ . We formalize this scenario – including the opposite case of high-frequency components being dominant – by studying  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)/\|\mathbf{F}(t)\|)$ , i.e. the Rayleigh quotient of  $\mathbf{I}_d \otimes \Delta$ .

**Definition 2.2.**  $\dot{\mathbf{F}}(t) = \text{GNN}_\theta(\mathbf{F}(t), t)$  initialized at  $\mathbf{F}(0)$  is *Low/High-Frequency-Dominant* (L/HFD) if  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)/\|\mathbf{F}(t)\|) \rightarrow 0$  (respectively,  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)/\|\mathbf{F}(t)\|) \rightarrow \rho_\Delta/2$ ) for  $t \rightarrow \infty$ .

We report a consequence of Definition 2.2 and refer to Appendix A.3 for additional details and motivations for the characterizations of LFD and HFD.

**Lemma 2.3.**  $\text{GNN}_\theta$  is LFD (HFD) iff for each  $t_j \rightarrow \infty$  there exist  $t_{j_k} \rightarrow \infty$  and  $\mathbf{F}_\infty$  s.t.  $\mathbf{F}(t_{j_k})/\|\mathbf{F}(t_{j_k})\| \rightarrow \mathbf{F}_\infty$  and  $\Delta \mathbf{f}_\infty^r = \mathbf{0}$  ( $\Delta \mathbf{f}_\infty^r = \rho_\Delta \mathbf{f}_\infty^r$ , respectively).

135 If a graph is *homophilic*, adjacent nodes are likely to share the same label and we expect a smoothing  
 136 or LFD dynamics enhancing the low-frequency components to be successful at node classification  
 137 tasks [43, 28]. In the opposite case of *heterophily*, the high-frequency components might contain more  
 138 relevant information for separating classes [4, 5] – the prototypical example being the eigenvector of  
 139  $\Delta$  associated with largest frequency  $\rho_\Delta$  separating a regular bipartite graph. In other words, the class  
 140 of heterophilic graphs contain instances where signals should be *sharpened* by increasing  $\mathcal{E}^{\text{Dir}}$  rather  
 141 than smoothed out. Accordingly, an ideal framework for learning on graphs must accommodate both  
 142 of these opposite scenarios by being able to induce either an LFD or a HFD dynamics.

143 **Parametric Dirichlet energy: channel-mixing as metric in feature space.** In eq. (1) a constant  
 144 nontrivial metric  $h$  in  $\mathbb{R}^d$  leads to the mixing of the feature channels. We adapt this idea by considering  
 145 a symmetric positive semi-definite  $\mathbf{H} = \mathbf{W}^\top \mathbf{W}$  with  $\mathbf{W} \in \mathbb{R}^{d \times d}$  and using it to generalize  $\mathcal{E}^{\text{Dir}}$  as

$$\mathcal{E}_{\mathbf{W}}^{\text{Dir}}(\mathbf{F}) := \frac{1}{4} \sum_{q,r=1}^d \sum_i \sum_{j:(i,j) \in \mathbf{E}} h_{qr}(\nabla \mathbf{f}^q)_{ij}(\nabla \mathbf{f}^r)_{ij} = \frac{1}{4} \sum_{(i,j) \in \mathbf{E}} \|\mathbf{W}(\nabla \mathbf{F})_{ij}\|^2. \quad (4)$$

146 We note the analogy with eq. (1), where the sum over the nodes replaces the integration over the  
 147 domain and the  $j$ -th derivative at some point  $i$  is replaced by the gradient along the edge  $(i, j) \in \mathbf{E}$ .  
 148 We generally treat  $\mathbf{W}$  as *learnable weights* and study the gradient flow of  $\mathcal{E}_{\mathbf{W}}^{\text{Dir}}$ :

$$\dot{\mathbf{F}}(t) = -\nabla_{\mathbf{F}} \mathcal{E}_{\mathbf{W}}^{\text{Dir}}(\mathbf{F}(t)) = -\Delta \mathbf{F}(t) \mathbf{W}^\top \mathbf{W}. \quad (5)$$

149 We see that eq. (5) generalizes eq. (3). Below ‘smoothing’ is intended as in Definition 2.1.

150 **Proposition 2.4.** Let  $P_{\mathbf{W}}^{\text{ker}}$  be the projection onto  $\ker(\mathbf{W}^\top \mathbf{W})$ . Equation (5) is smoothing since

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \leq e^{-2t \text{gap}(\mathbf{W}^\top \mathbf{W}) \text{gap}(\Delta)} \|\mathbf{F}(0)\|^2 + \mathcal{E}^{\text{Dir}}((P_{\mathbf{W}}^{\text{ker}} \otimes \mathbf{I}_n) \text{vec}(\mathbf{F}(0))), \quad t \geq 0.$$

151 In fact  $\mathbf{F}(t) \rightarrow \mathbf{F}_\infty$  s.t.  $\exists \phi_\infty \in \mathbb{R}^d$ : for each  $i \in \mathbf{V}$  we have  $(\mathbf{f}_\infty)_i = \sqrt{d_i} \phi_\infty + P_{\mathbf{W}}^{\text{ker}} \mathbf{f}_i(0)$ .

152 Proposition 2.4 implies that no weight matrix  $\mathbf{W}$  in eq. (5) can separate the limit embeddings  $\mathbf{F}(\infty)$   
 153 of nodes with same degree and input features. If  $\mathbf{W}$  has a trivial kernel, then nodes with same degrees  
 154 converge to the same representation and *over-smoothing* occurs as per Definition 2.1. Differently  
 155 from [29, 30, 8], over-smoothing occurs independently of the spectral radius of the ‘channel-mixing’  
 156 if its eigenvalues are *positive* – even for equations which lead to residual GNNs when discretized  
 157 [12]. According to Proposition 2.4, we do not expect eq. (5) to succeed on heterophilic graphs where  
 158 *smoothing* processes are generally harmful – this is confirmed in Figure 2 (see *prod*-curve). To  
 159 remedy this problem, we generalize eq. (5) to a gradient flow that can be HFD as per Definition 2.2.

### 160 3 A general parametric energy for pairwise interactions

161 We first rewrite the energy  $\mathcal{E}_{\mathbf{W}}^{\text{Dir}}$  in eq. (4) as

$$\mathcal{E}_{\mathbf{W}}^{\text{Dir}}(\mathbf{F}) = \frac{1}{2} \sum_i \langle \mathbf{f}_i, \mathbf{W}^\top \mathbf{W} \mathbf{f}_i \rangle - \frac{1}{2} \sum_{i,j} \bar{a}_{ij} \langle \mathbf{f}_i, \mathbf{W}^\top \mathbf{W} \mathbf{f}_j \rangle. \quad (6)$$

162 We then define a *new, more general* energy by replacing the occurrences of  $\mathbf{W}^\top \mathbf{W}$  with new  
 163 symmetric matrices  $\Omega, \mathbf{W} \in \mathbb{R}^{d \times d}$  since we also want to generate repulsive forces:

$$\mathcal{E}^{\text{tot}}(\mathbf{F}) := \frac{1}{2} \sum_i \langle \mathbf{f}_i, \Omega \mathbf{f}_i \rangle - \frac{1}{2} \sum_{i,j} \bar{a}_{ij} \langle \mathbf{f}_i, \mathbf{W} \mathbf{f}_j \rangle \equiv \mathcal{E}_\Omega^{\text{ext}}(\mathbf{F}) + \mathcal{E}_{\mathbf{W}}^{\text{pair}}(\mathbf{F}), \quad (7)$$

164 with associated gradient flow of the form (see Appendix B)

$$\dot{\mathbf{F}}(t) = -\nabla_{\mathbf{F}} \mathcal{E}^{\text{tot}}(\mathbf{F}(t)) = -\mathbf{F}(t) \Omega + \bar{\mathbf{A}} \mathbf{F}(t) \mathbf{W}. \quad (8)$$

165 Note that eq. (8) is gradient flow of some energy  $\mathbf{F} \mapsto \mathcal{E}^{\text{tot}}(\mathbf{F})$  iff both  $\Omega$  and  $\mathbf{W}$  are symmetric.

166 **A multi-particle system point of view: attraction vs repulsion.** Consider the  $d$ -dimensional  
 167 node-features as particles in  $\mathbb{R}^d$  with energy  $\mathcal{E}^{\text{tot}}$ . While the term  $\mathcal{E}_\Omega^{\text{ext}}$  is *independent of the graph*  
 168 *topology* and represents an **external** field in the feature space, the second term  $\mathcal{E}_{\mathbf{W}}^{\text{pair}}$  constitutes a  
 169 potential energy, with  $\mathbf{W}$  a *bilinear form* determining the **pairwise interactions** of adjacent node

representations. Given a symmetric  $\mathbf{W}$ , we write  $\mathbf{W} = \mathbf{\Theta}_+^\top \mathbf{\Theta}_+ - \mathbf{\Theta}_-^\top \mathbf{\Theta}_-$ , by decomposing the spectrum of  $\mathbf{W}$  in positive and negative values. We can rewrite  $\mathcal{E}^{\text{tot}} = \mathcal{E}_{\Omega-\mathbf{W}}^{\text{ext}} + \mathcal{E}_{\mathbf{\Theta}_+}^{\text{Dir}} - \mathcal{E}_{\mathbf{\Theta}_-}^{\text{Dir}}$ , i.e.

$$\mathcal{E}^{\text{tot}}(\mathbf{F}) = \frac{1}{2} \sum_i \langle \mathbf{f}_i, (\mathbf{\Omega} - \mathbf{W}) \mathbf{f}_i \rangle + \frac{1}{4} \sum_{i,j} \|\mathbf{\Theta}_+(\nabla \mathbf{F})_{ij}\|^2 - \frac{1}{4} \sum_{i,j} \|\mathbf{\Theta}_-(\nabla \mathbf{F})_{ij}\|^2. \quad (9)$$

The gradient flow of  $\mathcal{E}^{\text{tot}}$  minimizes  $\mathcal{E}_{\mathbf{\Theta}_+}^{\text{Dir}}$  and maximizes  $\mathcal{E}_{\mathbf{\Theta}_-}^{\text{Dir}}$ . The matrix  $\mathbf{W}$  encodes *repulsive pairwise interactions* via its negative-definite component  $\mathbf{\Theta}_-$  which lead to terms  $\|\mathbf{\Theta}_-(\nabla \mathbf{F})_{ij}\|$  increasing along the solution. The latter affords a ‘sharpening’ effect desirable on heterophilic graphs where we need to disentangle adjacent node representations and hence ‘magnify’ the edge-gradient.

**Spectral analysis of the channel-mixing.** We will now show that eq. (8) can lead to a HFD dynamics. To this end, we assume that  $\mathbf{\Omega} = \mathbf{0}$  so that eq. (8) becomes  $\dot{\mathbf{F}}(t) = \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$ . According to eq. (9) the negative eigenvalues of  $\mathbf{W}$  lead to repulsion. We show that the latter can induce HFD dynamics as per Definition 2.2. We let  $P_{\mathbf{W}}^{\rho_-}$  be the orthogonal projection into the eigenspace of  $\mathbf{W} \otimes \bar{\mathbf{A}}$  associated with the eigenvalue  $\rho_- := |\lambda_-^{\mathbf{W}}|(\rho_{\Delta} - 1)$ . We define  $\epsilon_{\text{HFD}}$  explicitly in eq. (24).

**Proposition 3.1.** *If  $\rho_- > \lambda_+^{\mathbf{W}}$ , then  $\dot{\mathbf{F}}(t) = \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$  is HFD for a.e.  $\mathbf{F}(0)$ : there exists  $\epsilon_{\text{HFD}}$  s.t.*

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) = e^{2t\rho_-} \left( \frac{\rho_{\Delta}}{2} \|\mathbf{P}_{\mathbf{W}}^{\rho_-} \mathbf{F}(0)\|^2 + \mathcal{O}(e^{-2t\epsilon_{\text{HFD}}}) \right), \quad t \geq 0,$$

and  $\mathbf{F}(t)/\|\mathbf{F}(t)\|$  converges to  $\mathbf{F}_{\infty} \in \mathbb{R}^{n \times d}$  such that  $\Delta \mathbf{f}_{\infty}^r = \rho_{\Delta} \mathbf{f}_{\infty}^r$ , for  $1 \leq r \leq d$ .

Proposition 3.1 shows that *if enough mass of the spectrum of the ‘channel-mixing’ is distributed over the negative eigenvalues, then the evolution is dominated by the graph high frequencies*. This analysis is made possible in our gradient flow framework where  $\mathbf{W}$  must be *symmetric*. The HFD dynamics induced by negative eigenvalues of  $\mathbf{W}$  is confirmed in Figure 2 (*neg-prod-curve* in the bottom chart).

**A more general energy.** Equations with a source term may have better expressive power [44, 11, 39]. In our framework this means adding an extra energy term of the form  $\mathcal{E}_{\tilde{\mathbf{W}}}^{\text{source}}(\mathbf{F}) := \beta \langle \mathbf{F}, \mathbf{F}(0) \tilde{\mathbf{W}} \rangle$  to eq. (7) with some learnable  $\beta$  and  $\tilde{\mathbf{W}}$ . This leads to the following gradient flow:

$$\dot{\mathbf{F}}(t) = -\mathbf{F}(t)\mathbf{\Omega} + \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W} - \beta\mathbf{F}(0)\tilde{\mathbf{W}}. \quad (10)$$

We also observe that one could replace the fixed matrix  $\bar{\mathbf{A}}$  with a more general *symmetric graph vector field*  $\mathcal{A}$  satisfying  $\mathcal{A}_{ij} = 0$  if  $(i, j) \notin \mathbf{E}$ , although in this work we focus on the case  $\mathcal{A} = \bar{\mathbf{A}}$ . We also note that when  $\mathbf{\Omega} = \mathbf{W}$ , then eq. (8) becomes  $\dot{\mathbf{F}}(t) = -\Delta \mathbf{F}(t)\mathbf{W}$ . We perform a spectral analysis of this case in Appendix B.2.

**Non-linear activations.** In Appendix B.3 we discuss non-linear gradient flow equations. Here we study what happens if the gradient flow in eq. (10) is activated *pointwise* by  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ . We show that although we are no longer a gradient flow, the learnable multi-particle energy  $\mathcal{E}^{\text{tot}}$  is still decreasing along the solution, meaning that the interpretation of the channel-mixing  $\mathbf{W}$  inducing attraction and repulsion via its positive and negative eigenvalues respectively is **preserved**.

**Proposition 3.2.** *Consider a non-linear map  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  such that the function  $x \mapsto x\sigma(x) \geq 0$ . If  $t \mapsto \mathbf{F}(t)$  solves the equation*

$$\dot{\mathbf{F}}(t) = \sigma \left( -\mathbf{F}(t)\mathbf{\Omega} + \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W} - \beta\mathbf{F}(0)\tilde{\mathbf{W}} \right),$$

where  $\sigma$  acts elementwise, then

$$\frac{d\mathcal{E}^{\text{tot}}(\mathbf{F}(t))}{dt} \leq 0.$$

A proof of this result and more details and discussion are reported in Appendix E. We emphasize here that differently from previous results about behaviour of ReLU wrt  $\mathcal{E}^{\text{Dir}}$  [30, 8], we deal with a much more general energy that can also induce repulsion and a more general family of activation functions (that include ReLU, tanh, arctan and many others).

## 4 Comparison with GNNs

In this Section, we study standard GNN models from the perspective of our gradient flow framework.

## 208 4.1 Continuous case

209 Continuous GNN models replace layers with continuous time. In contrast with Proposition 3.1,  
 210 we show that three main *linearized* continuous GNN models are either *smoothing* or LFD as  
 211 per Definition 2.2. The linearized PDE-GCN<sub>D</sub> model [17] corresponds to choosing  $\beta = 0$  and  
 212  $\Omega = \mathbf{W} = \mathbf{K}(t)^\top \mathbf{K}(t)$  in eq. (10), for some time-dependent family  $t \mapsto \mathbf{K}(t) \in \mathbb{R}^{d \times d}$ :

$$\dot{\mathbf{F}}_{\text{PDE-GCN}_D}(t) = -\Delta \mathbf{F}(t) \mathbf{K}(t)^\top \mathbf{K}(t).$$

213 The CGNN model [44] can be derived from eq. (10) by setting  $\Omega = \mathbf{I} - \tilde{\Omega}$ ,  $\mathbf{W} = \tilde{\mathbf{W}} = \mathbf{I}$ ,  $\beta = 1$ :

$$\dot{\mathbf{F}}_{\text{CGNN}}(t) = -\Delta \mathbf{F}(t) + \mathbf{F}(t) \tilde{\Omega} + \mathbf{F}(0).$$

214 Finally, in linearized GRAND [10] a row-stochastic matrix  $\mathcal{A}(\mathbf{F}(0))$  is *learned* from the encoding  
 215 via an attention mechanism and we have

$$\dot{\mathbf{F}}_{\text{GRAND}}(t) = -\Delta_{\text{RW}} \mathbf{F}(t) = -(\mathbf{I} - \mathcal{A}(\mathbf{F}(0))) \mathbf{F}(t).$$

216 We note that if  $\mathcal{A}$  is not symmetric, then GRAND is *not* a gradient flow.

217 **Proposition 4.1.** PDE – GCN<sub>D</sub>, CGNN and GRAND satisfy the following:

- 218 (i) PDE – GCN<sub>D</sub> is a *smoothing model*:  $\dot{\mathcal{E}}^{\text{Dir}}(\mathbf{F}_{\text{PDE-GCN}_D}(t)) \leq 0$ .
- 219 (ii) For a.e.  $\mathbf{F}(0)$  it holds: CGNN is never HFD and if we remove the source term, then  
 220  $\mathcal{E}^{\text{Dir}}(\mathbf{F}_{\text{CGNN}}(t)/\|\mathbf{F}_{\text{CGNN}}(t)\|) \leq e^{-\text{gap}(\Delta)t}$ .
- 221 (iii) If  $G$  is connected,  $\mathbf{F}_{\text{GRAND}}(t) \rightarrow \boldsymbol{\mu}$  as  $t \rightarrow \infty$ , with  $\boldsymbol{\mu}^r = \text{mean}(\mathbf{f}^r(0))$ ,  $1 \leq r \leq d$ .

222 By (ii) the source-free CGNN-evolution is LFD *independent of*  $\tilde{\Omega}$ . Moreover, by (iii), over-smoothing  
 223 occurs for GRAND as per Definition 2.1. On the other hand, Proposition 3.1 shows that the negative  
 224 eigenvalues of  $\mathbf{W}$  can make the source-free gradient flow in eq. (8) HFD. Experiments in Section 5  
 225 confirm that the gradient flow model outperforms CGNN and GRAND on heterophilic graphs.

## 226 4.2 Discrete case

227 We now describe a discrete version of our gradient flow model and compare it to ‘discrete’ GNNs  
 228 where discrete time steps correspond to different layers. In the spirit of [12], we use explicit Euler  
 229 scheme with step size  $\tau \leq 1$  to solve eq. (10) and set  $\tilde{\mathbf{W}} = \mathbf{I}$ . In the gradient flow framework we  
 230 *parametrize the energy* rather than the actual equations, which leads to *symmetric* channel-mixing  
 231 matrices  $\Omega, \mathbf{W} \in \mathbb{R}^{d \times d}$  that are *shared across the layers*. Since the matrices are square, an *encoding*  
 232 block  $\psi_{\text{EN}} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times d}$  is used to process input features  $\mathbf{F}_0 \in \mathbb{R}^{n \times p}$  and generally reduce the  
 233 hidden dimension from  $p$  to  $d$ . Moreover, the iterations inherently lead to a residual architecture  
 234 because of the explicit Euler discretization:

$$\mathbf{F}(t + \tau) = \mathbf{F}(t) + \tau (-\mathbf{F}(t) \Omega + \bar{\mathbf{A}} \mathbf{F}(t) \mathbf{W} + \beta \mathbf{F}(0)), \quad \mathbf{F}(0) = \psi_{\text{EN}}(\mathbf{F}_0), \quad (11)$$

235 with prediction  $y = \psi_{\text{DE}}(\mathbf{F}(T))$  produced by a *decoder*  $\psi_{\text{DE}} : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times k}$ , where  $k$  is the  
 236 number of label classes and  $T$  *integration time* of the form  $T = m\tau$ , so that  $m \in \mathbb{N}$  represents the  
 237 number of *layers*. Although eq. (11) is linear, we can include non-linear activations in  $\psi_{\text{EN}}, \psi_{\text{DE}}$   
 238 making the entire model generally non-linear. We emphasize two important points:

- 239 • Since the framework is residual, even if the message-passing is linear, this is *not equivalent*  
 240 to collapsing the dynamics into a single layer with diffusion matrix  $\bar{\mathbf{A}}^m$ , with  $m$  the number  
 241 of layers, see eq. (27) in the appendix where we derive the expansion of the solution.
- 242 • We could also activate the equations pointwise and maintain the physics interpretation thanks  
 243 to Proposition 3.2 to gain greater expressive power. In the following though, we mainly  
 244 stick to the linear discrete gradient flow unless otherwise stated.

245 **Are discrete GNNs gradient flows?** Given a (learned) symmetric graph vector field  $\mathcal{A} \in \mathbb{R}^{n \times n}$   
 246 satisfying  $\mathcal{A}_{ij} = 0$  if  $(i, j) \notin E$ , consider a family of linear GNNs with shared weights of the form

$$\mathbf{F}(t + 1) = \mathbf{F}(t) \Omega + \mathcal{A} \mathbf{F}(t) \mathbf{W} + \beta \mathbf{F}(0) \tilde{\mathbf{W}}, \quad 0 \leq t \leq T. \quad (12)$$

247 Symmetry is the key requirement to interpret GNNs in eq. (12) in a gradient flow framework.

**Lemma 4.2.** Equation (12) is the unit step size discrete gradient flow of  $\mathcal{E}_{\mathbf{I}-\mathbf{\Omega}}^{\text{ext}} + \mathcal{E}_{\mathbf{A},\mathbf{W}}^{\text{pair}} - \mathcal{E}_{\mathbf{W}}^{\text{source}}$ , with  $\mathcal{E}_{\mathbf{A},\mathbf{W}}^{\text{pair}}$  defined by replacing  $\bar{\mathbf{A}}$  with  $\mathbf{A}$  in eq. (7), iff  $\mathbf{\Omega}$  and  $\mathbf{W}$  are symmetric.

Lemma 4.2 provides a recipe for making standard architectures into a gradient flow, with *symmetry* being the key requirement. When eq. (12) is a gradient flow, the underlying GNN dynamics is equivalent to minimizing a multi-particle energy by learning attractive and repulsive directions in feature space as discussed in Section 3. In Appendix C.2, we show how Lemma 4.2 covers linear versions of GCN [27, 43], GAT [42], GraphSAGE [23] and GCNII [11] to name a few.

**Over-smoothing analysis in discrete setting.** By Proposition 3.1 we know that the continuous version of eq. (11) can be HFD thanks to the negative eigenvalues of  $\mathbf{W}$ . The next result represents a discrete counterpart of Proposition 3.1 and shows that *residual, symmetrized graph convolutional models can be HFD*. Below  $P_{\mathbf{W}}^{\rho_-}$  is the projection into the eigenspace associated with the eigenvalue  $\rho_- := |\lambda_-^{\mathbf{W}}|(\rho_{\Delta} - 1)$  and we report the explicit value of  $\delta_{\text{HFD}}$  in eq. (28) in Appendix C.3. We let:

$$\lambda_+^{\mathbf{W}}(\rho_{\Delta} - 1)^{-1} < |\lambda_-^{\mathbf{W}}| < 2(\tau(2 - \rho_{\Delta}))^{-1}. \quad (13)$$

**Theorem 4.3.** Given  $\mathbf{F}(t + \tau) = \mathbf{F}(t) + \tau \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$ , with  $\mathbf{W}$  symmetric, if eq. (13) holds then

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}(m\tau)) = (1 + \tau\rho_-)^{2m} \left( \frac{\rho_{\Delta}}{2} \|P_{\mathbf{W}}^{\rho_-} \mathbf{F}(0)\|^2 + \mathcal{O} \left( \left( \frac{1 + \tau\delta_{\text{HFD}}}{1 + \tau\rho_-} \right)^{2m} \right) \right), \quad \delta_{\text{HFD}} < \rho_-,$$

hence the dynamics is HFD for a.e.  $\mathbf{F}(0)$  and in fact  $\mathbf{F}(m\tau)/\|\mathbf{F}(m\tau)\| \rightarrow \mathbf{F}_{\infty}$  s.t.  $\Delta \mathbf{f}_{\infty}^r = \rho_{\Delta} \mathbf{f}_{\infty}^r$ . Conversely, if  $\mathbf{G}$  is not bipartite, then for a.e.  $\mathbf{F}(0)$  the system  $\mathbf{F}(t + \tau) = \tau \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$ , with  $\mathbf{W}$  symmetric, is LFD independent of the spectrum of  $\mathbf{W}$ .

Theorem 4.3 shows that linear discrete gradient flows can be HFD due to the negative eigenvalues of  $\mathbf{W}$ . This differs from statements that standard GCNs act as low-pass filters and thus over-smooth in the limit. Indeed, in these cases the spectrum of  $\mathbf{W}$  is generally ignored [43, 11] or required to be sufficiently small in terms of singular value decomposition [29, 30, 8] *when no residual connection is present*. On the other hand, Theorem 4.3 emphasizes that the spectrum of  $\mathbf{W}$  plays a key role to enhance the high frequencies when enough mass is distributed over the negative eigenvalues provided that a residual connection exists – this is confirmed by the *neg-prod*-curve in Figure 2.

**The residual connection from a spectral perspective.** Given a sufficiently small step-size so that the right hand side of inequality 13 is satisfied,  $\mathbf{F}(t + \tau) = \mathbf{F}(t) + \tau \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$  is HFD for a.e.  $\mathbf{F}(0)$  if  $|\lambda_-^{\mathbf{W}}|(\rho_{\Delta} - 1) > \lambda_+^{\mathbf{W}}$ , i.e. ‘there is more mass’ in the negative spectrum of  $\mathbf{W}$  than in the positive one. This means that differently from [29, 30, 8], there is no requirement on the minimal magnitude of the spectral radius of  $\mathbf{W}$  coming from the graph topology as long as  $\lambda_+^{\mathbf{W}}$  is small enough. Conversely, without a residual term, the dynamics is LFD for a.e.  $\mathbf{F}(0)$  *independently* of the sign and magnitude of the eigenvalues of  $\mathbf{W}$ . This is also confirmed by the GCN-curve in Figure 2.

**Over-smoothing vs LFD.** We highlight how in general a linear GCN equation as  $\mathbf{F}(t + \tau) = \tau \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$  may avoid over-smoothing in the sense of Definition 2.1, meaning that  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow \infty$  as soon as there exist  $\lambda_i^{\Delta} \in (0, 1)$  and the spectral radius of  $\mathbf{W}$  is large enough. However, this will not lead to over-separation since the dominating term is the lowest frequency one: in other words, once we re-set the scale right as per the normalization in Theorem 4.3, we encounter loss of separability even with large (and possibly negative) spectrum of  $\mathbf{W}$ .

## 5 Experiments

In this section we evaluate the gradient flow framework (GRAFF). We corroborate the spectral analysis using synthetic data with controllable homophily. We confirm that having negative (positive) eigenvalues of the channel-mixing  $\mathbf{W}$  are essential in heterophilic (homophilic) scenarios where the gradient flow should align with HFD (LFD) respectively. We show that the gradient flow in eq. (11) – a linear, residual, symmetric graph convolutional model – achieves competitive performance on heterophilic datasets.

**Methodology.** We crystallize GRAFF in the model presented in eq. (11) with  $\psi_{\text{EN}}, \psi_{\text{DE}}$  implemented as single linear layers or MLPs, and we set  $\Omega$  to be diagonal. For the real-world experiments we consider *diagonally-dominant* (DD), *diagonal* (D) and *time-dependent* choices for the structure of  $\mathbf{W}$  that offer explicit control over its spectrum. In the (DD)-case, we consider a  $\mathbf{W}^0 \in \mathbb{R}^{d \times d}$  symmetric with zero diagonal and  $\mathbf{w} \in \mathbb{R}^d$  defined by  $\mathbf{w}_\alpha = q_\alpha \sum_\beta |\mathbf{W}_{\alpha\beta}^0| + r_\alpha$ , and set  $\mathbf{W} = \text{diag}(\mathbf{w}) + \mathbf{W}^0$ . Due to the Gershgorin Theorem the eigenvalues of  $\mathbf{W}$  belong to  $[\mathbf{w}_\alpha - \sum_\beta |\mathbf{W}_{\alpha\beta}^0|, \mathbf{w}_\alpha + \sum_\beta |\mathbf{W}_{\alpha\beta}^0|]$ , so the model ‘can’ easily re-distribute mass in the spectrum of  $\mathbf{W}$  via  $q_\alpha, r_\alpha$ . This generalizes the decomposition of  $\mathbf{W}$  in [11] providing a justification in terms of its spectrum and turns out to be more efficient w.r.t. the hidden dimension  $d$  as shown in Figure 4 in the Appendix. For (D) we take  $\mathbf{W}$  to be diagonal, with entries sampled  $\mathcal{U}[-1, 1]$  and fixed – i.e., **we do not train** over  $\mathbf{W}$  – and only learn  $\psi_{\text{EN}}, \psi_{\text{DE}}$ . We also include a *time-dependent* model where  $\mathbf{W}_t$  varies across layers. To investigate the role of the spectrum of  $\mathbf{W}$  on synthetic graphs, we construct three additional variants:  $\mathbf{W} = \mathbf{W}' + \mathbf{W}'^\top$ ,  $\mathbf{W} = \pm \mathbf{W}'^\top \mathbf{W}'$  named *sum*, *prod* and *neg-prod* respectively where *prod* (*neg-prod*) variants have only non-negative (non-positive) eigenvalues.

**Complexity and number of parameters.** If we treat the number of layers as a constant, the discrete gradient flow scales as  $\mathcal{O}(|V|pd + |E|d^2)$ , where  $p$  and  $d$  are input feature and hidden dimension respectively, with  $p \geq d$  usually. Note that GCN has complexity  $\mathcal{O}(|E|pd)$  and in fact *our model is faster than GCN* as confirmed in Figure 5 in Appendix D. Since  $\psi_{\text{EN}}, \psi_{\text{DE}}$  are single linear layers (MLPs), we can bound the number of parameters by  $pd + d^2 + 3d + dk$ , with  $k$  the number of label classes, in the (DD)-variant while in the (D)-variant we have  $pd + 3d + dk$ . Further ablation studies appear in Figure 4 in the Appendix showing that (DD) outperforms *sum* and GCN – especially in the lower hidden dimension regime – on real-world benchmarks with varying homophily.

### Synthetic experiments and ablation studies.

To investigate our claims in a controlled environment we use the synthetic Cora dataset of [51, Appendix G]. Graphs are generated for target levels of homophily via preferential attachment – see Appendix D.3 for details. Figure 2 confirms the spectral analysis and offers a better understanding in terms of performance and smoothness of the predictions. Each curve – except GCN – represents one version of  $\mathbf{W}$  as in ‘methodology’ and we implement eq. (11) with  $\beta = 0, \Omega = \mathbf{0}$ . Figure 2 (top) reports the test accuracy vs true label homophily. *Neg-prod* is better than *prod* on low-homophily and viceversa on high-homophily. This confirms Proposition 3.1 where we have shown that the gradient flow can lead to a HFD dynamics – that are generally desirable with low-homophily – through the negative eigenvalues of  $\mathbf{W}$ . Conversely, the *prod* configuration (where we have an attraction-only dynamics) struggles in low-homophily scenarios *even though a residual connection is present*. Both *prod* and *neg-prod* are ‘extreme’ choices and serve the purpose of highlighting that by turning off one side of the spectrum this could be the more damaging depending on the underlying homophily. In general though ‘neutral’ variants like *sum* and (DD) are indeed more flexible and better performing. In fact, (DD) outperforms GCN especially in low-homophily scenarios, confirming Theorem 4.3 where we have shown that without a residual connection convolutional models are LFD – and hence more sensitive to underlying homophily – irrespectively of the spectrum of  $\mathbf{W}$ . This is further confirmed in Figure 3.

In Figure 2 (bottom) we compute the homophily of the prediction (cross) for a given method and we compare with the homophily (circle) of the prediction read from the encoding (i.e. *graph-agnostic*). The homophily here is a proxy to assess whether the evolution is *smoothing*, the goal being explaining the smoothness of the prediction via the spectrum of  $\mathbf{W}$  as per our theoretical analysis. For *neg-prod* the homophily after the evolution is lower than that of the encoding, supporting the analysis that negative eigenvalues of  $\mathbf{W}$  enhance high-frequencies. The opposite behaviour occurs in the case of *prod* and explains that in the low-homophily regime *prod* is under-performant due to the prediction

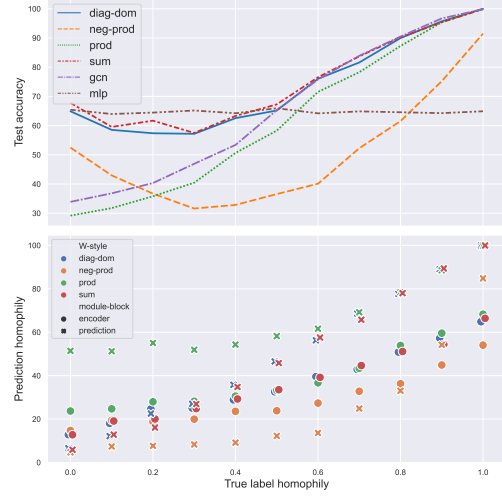


Figure 2: Experiments on synthetic datasets with controlled homophily.



	Texas 0.11	Wisconsin 0.21	Cornell 0.30	Film 0.22	Squirrel 0.22	Chameleon 0.23	Citeseer 0.74	Pubmed 0.80	Cora 0.81
Hom level	183	251	183	7,600	5,201	2,277	3,327	18,717	2,708
#Nodes	295	466	280	26,752	198,493	31,421	4,676	44,327	5,278
#Edges	5	5	5	5	5	5	7	3	6
#Classes									
GGCN	84.86 ± 4.55	86.86 ± 3.29	85.68 ± 6.63	37.54 ± 1.56	55.17 ± 1.58	71.14 ± 1.84	77.14 ± 1.45	89.15 ± 0.37	87.95 ± 1.05
GPRGNN	78.38 ± 4.36	82.94 ± 4.21	80.27 ± 8.11	34.63 ± 1.22	31.61 ± 1.24	46.58 ± 1.71	77.13 ± 1.67	87.54 ± 0.38	87.95 ± 1.18
H2GCN	84.86 ± 7.23	87.65 ± 4.98	82.70 ± 5.28	35.70 ± 1.00	36.48 ± 1.86	60.11 ± 2.15	77.11 ± 1.57	89.49 ± 0.38	87.87 ± 1.20
GCNII	77.57 ± 3.83	80.39 ± 3.40	77.86 ± 3.79	37.44 ± 1.30	38.47 ± 1.58	63.86 ± 3.04	77.33 ± 1.48	90.15 ± 0.43	88.37 ± 1.25
Geom-GCN	66.76 ± 2.72	64.51 ± 3.66	60.54 ± 3.67	31.59 ± 1.15	38.15 ± 0.92	60.00 ± 2.81	78.02 ± 1.15	89.95 ± 0.47	85.35 ± 1.57
PairNorm	60.27 ± 4.34	48.43 ± 6.14	58.92 ± 3.15	27.40 ± 1.24	50.44 ± 2.04	62.74 ± 2.82	73.59 ± 1.47	87.53 ± 0.44	85.79 ± 1.01
GraphSAGE	82.43 ± 6.14	81.18 ± 5.56	75.95 ± 5.01	34.23 ± 0.99	41.61 ± 0.74	58.73 ± 1.68	76.04 ± 1.30	88.45 ± 0.50	86.90 ± 1.04
GCN	55.14 ± 5.16	51.76 ± 3.06	60.54 ± 5.30	27.32 ± 1.10	53.43 ± 2.01	64.82 ± 2.24	76.50 ± 1.36	88.42 ± 0.50	86.98 ± 1.27
GAT	52.16 ± 6.63	49.41 ± 4.09	61.89 ± 5.05	27.44 ± 0.89	40.72 ± 1.55	60.26 ± 2.50	76.55 ± 1.23	87.30 ± 1.10	86.33 ± 0.48
MLP	80.81 ± 4.75	85.29 ± 3.31	81.89 ± 6.40	36.53 ± 0.70	28.77 ± 1.56	46.21 ± 2.99	74.02 ± 1.90	75.69 ± 2.00	87.16 ± 0.37
CGNN	71.35 ± 4.05	74.31 ± 7.26	66.22 ± 7.69	35.95 ± 0.86	29.24 ± 1.09	46.89 ± 1.66	76.91 ± 1.81	87.70 ± 0.49	87.10 ± 1.35
GRAND	75.68 ± 7.25	79.41 ± 3.64	82.16 ± 7.09	35.62 ± 1.01	40.05 ± 1.50	54.67 ± 2.54	76.46 ± 1.77	89.02 ± 0.51	87.36 ± 0.96
Sheaf (max)	85.95 ± 5.51	89.41 ± 4.74	84.86 ± 4.71	37.81 ± 1.15	56.34 ± 1.32	68.04 ± 1.58	76.70 ± 1.57	89.49 ± 0.40	86.90 ± 1.13
GRAFF (DD)	88.38 ± 4.53	87.45 ± 2.94	83.24 ± 6.49	36.09 ± 0.81	54.52 ± 1.37	71.08 ± 1.75	76.92 ± 1.70	88.95 ± 0.52	87.61 ± 0.97
GRAFF (D)	88.11 ± 5.57	88.83 ± 3.29	84.05 ± 6.10	37.11 ± 1.08	47.36 ± 1.89	66.78 ± 1.28	77.30 ± 1.85	90.04 ± 0.41	88.01 ± 1.03
GRAFF-timedep (DD)	87.03 ± 4.49	87.06 ± 4.04	82.16 ± 7.07	35.93 ± 1.23	53.97 ± 1.45	69.56 ± 1.20	76.59 ± 1.53	88.26 ± 0.41	87.38 ± 1.05

Table 1: Results on heterophilic and homophilic datasets

being smoother than the true homophily. (DD) and *sum* variants adapt better to the true homophily. We note how the encoding compensates when the dynamics can only either attract or repulse (i.e. the spectrum of  $\mathbf{W}$  has a sign) by decreasing or increasing the initial homophily respectively.

**Real world experiments.** We test GRAFF against a range of datasets with varying homophily [37, 33, 31] (see Appendix D.4 for additional details). We use results provided in [45, Table 1], which includes standard baselines as GCN [27], GraphSAGE [23], GAT [42], PairNorm [48] and recent models tailored towards the heterophilic setting (GGCN [45], Geom-GCN [31], H2GCN [51] and GPRGNN [13]). For Sheaf [5], a recent top-performer on heterophilic datasets, we took the best performing variant (out of six provided) for each dataset. We also include continuous baselines CGNN [44] and GRAND [10] to provide empirical evidence for Proposition 4.1. Splits taken from [31] are used in all the comparisons. The GRAFF model discussed in ‘methodology’ is a very simple architecture with shared parameters across layers and run-time smaller than GCN and more recent models like GGCN designed for heterophilic graphs (see Figure 5 in the Appendix). Nevertheless, it achieves competitive results on all datasets, performing on par or better than more complex recent models. Moreover, comparison with the ‘time-dependent’ (DD) variant confirms that by sharing weights across layers we do not lose performance. We note that on heterophilic graphs short integration time is usually needed due to the topology being harmful and the negative eigenvalues of  $\mathbf{W}$  leading to exponential behaviour (see Appendix D).

## 6 Conclusions

In this work, we developed a framework for GNNs where the evolution can be interpreted as minimizing a multi-particle learnable energy. This translates into studying the interaction between the spectrum of the graph and the spectrum of the ‘channel-mixing’ leading to a better understanding of when and why the induced dynamics is low (high) frequency dominated. From a theoretical perspective, we refined existing asymptotic analysis of GNNs to account for the role of the spectrum of the channel-mixing as well. From a practical perspective, our framework allows for ‘educated’ choices resulting in a simple convolutional model that achieves competitive performance on homophilic and heterophilic benchmarks while being faster than GCN. Our results refute the folklore of graph convolutional models being too simple for heterophilic benchmarks.

**Limitations and future works.** We limited our attention to a *constant* bilinear form  $\mathbf{W}$ , which might be excessively rigid. It is possible to derive non-constant alternatives that are *aware* of the features or the position in the graph. The main challenge amounts to matching the requirement for local ‘heterogeneity’ with efficiency: we reserve this question for future work. Our analysis is also a first step into studying the interaction of the graph and ‘channel-mixing’ spectra; we did not explore other dynamics that are neither LFD nor HFD as per our definitions. The energy formulation points to new models more ‘physics’ inspired; this will be explored in future work.

**Societal impact.** Our work sheds light on the actual dynamics of GNNs and could hence improve their understanding, which is crucial for assessing their impact on large-scale applications. We also show that instances of our framework achieve competitive performance on heterophilic data despite being faster than GCN, providing evidence for efficient methods with reduced footprint.

## References

- [1] U. Alon and E. Yahav. On the bottleneck of graph neural networks and its practical implications. In *International Conference on Learning Representations*, 2021.
- [2] M. Balcilar, G. Renton, P. Héroux, B. Gaüzère, S. Adam, and P. Honeine. Analyzing the expressive power of graph neural networks in a spectral perspective. In *International Conference on Learning Representations*, 2020.
- [3] M. Biloš, J. Sommer, S. S. Rangapuram, T. Januschowski, and S. Günnemann. Neural flows: Efficient alternative to neural odes. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [4] D. Bo, X. Wang, C. Shi, and H. Shen. Beyond low-frequency information in graph convolutional networks. In *AAAI AAAI Press*, 2021.
- [5] C. Bodnar, F. Di Giovanni, B. P. Chamberlain, P. Liò, and M. M. Bronstein. Neural sheaf diffusion: A topological perspective on heterophily and oversmoothing in gnns. *arXiv preprint arXiv:2202.04579*, 2022.
- [6] S. Brody, U. Alon, and E. Yahav. How attentive are graph attention networks? *arXiv preprint arXiv:2105.14491*, 2021.
- [7] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. In *2nd International Conference on Learning Representations, ICLR 2014*, 2014.
- [8] C. Cai and Y. Wang. A note on over-smoothing for graph neural networks. *arXiv preprint arXiv:2006.13318*, 2020.
- [9] B. Chamberlain, J. Rowbottom, D. Eynard, F. Di Giovanni, X. Dong, and M. Bronstein. Beltrami flow and neural diffusion on graphs. *Advances in Neural Information Processing Systems*, 34, 2021.
- [10] B. Chamberlain, J. Rowbottom, M. I. Gorinova, M. Bronstein, S. Webb, and E. Rossi. Grand: Graph neural diffusion. In *International Conference on Machine Learning*, pages 1407–1418. PMLR, 2021.
- [11] M. Chen, Z. Wei, Z. Huang, B. Ding, and Y. Li. Simple and deep graph convolutional networks. In *International Conference on Machine Learning*, pages 1725–1735. PMLR, 2020.
- [12] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- [13] E. Chien, J. Peng, P. Li, and O. Milenkovic. Adaptive universal generalized pagerank graph neural network. In *9th International Conference on Learning Representations, ICLR 2021*, 2021.
- [14] F. R. Chung and F. C. Graham. *Spectral graph theory*. Number 92. American Mathematical Soc., 1997.
- [15] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29, 2016.
- [16] J. Eells and J. H. Sampson. Harmonic mappings of riemannian manifolds. *American journal of mathematics*, 86(1):109–160, 1964.
- [17] M. Eliasof, E. Haber, and E. Treister. Pde-gcn: Novel architectures for graph neural networks motivated by partial differential equations. *Advances in Neural Information Processing Systems*, 34, 2021.
- [18] M. Geiger, L. Petrini, and M. Wyart. Landscape and training regimes in deep learning. *Physics Reports*, 924:1–18, 2021.

- [19] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning*, pages 1263–1272. PMLR, 2017.
- [20] C. Goller and A. Kuchler. Learning task-dependent distributed representations by backpropagation through structure. In *Proceedings of International Conference on Neural Networks (ICNN’96)*, volume 1, pages 347–352. IEEE, 1996.
- [21] M. Gori, G. Monfardini, and F. Scarselli. A new model for learning in graph domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 2, pages 729–734. IEEE, 2005.
- [22] E. Haber and L. Ruthotto. Stable architectures for deep neural networks. *Inverse Problems*, 34, 2018.
- [23] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.
- [24] D. K. Hammond, P. Vandergheynst, and R. Gribonval. The spectral graph wavelet transform: Fundamental theory and fast computation. In *Vertex-Frequency Analysis of Graph Signals*, pages 141–175. Springer, 2019.
- [25] M. He, Z. Wei, H. Xu, et al. Bernnet: Learning arbitrary graph spectral filters via bernstein approximation. *Advances in Neural Information Processing Systems*, 34, 2021.
- [26] R. Kimmel, N. Sochen, and R. Malladi. From high energy physics to low level vision. In *International Conference on Scale-Space Theories in Computer Vision*, pages 236–247. Springer, 1997.
- [27] T. N. Kipf and M. Welling. Semi-Supervised Classification with Graph Convolutional Networks. In *Proceedings of the 5th International Conference on Learning Representations, ICLR ’17*, 2017.
- [28] J. Klicpera, S. Weißenberger, and S. Günnemann. Diffusion improves graph learning. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019.
- [29] H. Nt and T. Maehara. Revisiting graph neural networks: All we have is low-pass filters. *arXiv preprint arXiv:1905.09550*, 2019.
- [30] K. Oono and T. Suzuki. Graph neural networks exponentially lose expressive power for node classification. In *International Conference on Learning Representations*, 2020.
- [31] H. Pei, B. Wei, K. C. Chang, Y. Lei, and B. Yang. Geom-gcn: Geometric graph convolutional networks. In *8th International Conference on Learning Representations, ICLR 2020*, 2020.
- [32] P. Perona and J. Malik. Scale-space and edge detection using anisotropic diffusion. *PAMI*, 12(7):629–639, 1990.
- [33] B. Rozemberczki, C. Allen, and R. Sarkar. Multi-scale attributed node embedding. *Journal of Complex Networks*, 9(2):cnab014, 2021.
- [34] T. K. Rusch, B. P. Chamberlain, J. Rowbottom, S. Mishra, and M. M. Bronstein. Graph-coupled oscillator networks. In *International Conference on Machine Learning*, 2022.
- [35] M. E. Sander, P. Ablin, M. Blondel, and G. Peyré. Sinkformers: Transformers with doubly stochastic attention. In *International Conference on Artificial Intelligence and Statistics*, pages 3515–3530. PMLR, 2022.
- [36] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
- [37] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.

- [38] A. Sperduti. Encoding labeled graphs by labeling raam. *Advances in Neural Information Processing Systems*, 6, 1993.
- [39] M. Thorpe, T. M. Nguyen, H. Xia, T. Strohmer, A. Bertozzi, S. Osher, and B. Wang. Grand++: Graph neural diffusion with a source term. In *International Conference on Learning Representations*, 2021.
- [40] J. Topping, F. Di Giovanni, B. P. Chamberlain, X. Dong, and M. M. Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature. *International Conference on Learning Representations*, 2022.
- [41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [42] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- [43] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. PMLR, 2019.
- [44] L.-P. Xhonneux, M. Qu, and J. Tang. Continuous graph neural networks. In *International Conference on Machine Learning*, pages 10432–10441. PMLR, 2020.
- [45] Y. Yan, M. Hashemi, K. Swersky, Y. Yang, and D. Koutra. Two sides of the same coin: Heterophily and oversmoothing in graph convolutional neural networks. *arXiv preprint arXiv:2102.06462*, 2021.
- [46] Z. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec. Gnnexplainer: Generating explanations for graph neural networks. *Advances in neural information processing systems*, 32, 2019.
- [47] H. Yuan, H. Yu, S. Gui, and S. Ji. Explainability in graph neural networks: A taxonomic survey. *arXiv preprint arXiv:2012.15445*, 2020.
- [48] L. Zhao and L. Akoglu. Pairnorm: Tackling oversmoothing in gnns. *arXiv preprint arXiv:1909.12223*, 2019.
- [49] D. Zhou and B. Schölkopf. Regularization on discrete spaces. In *Joint Pattern Recognition Symposium*, pages 361–368. Springer, 2005.
- [50] K. Zhou, X. Huang, D. Zha, R. Chen, L. Li, S.-H. Choi, and X. Hu. Dirichlet energy constrained learning for deep graph neural networks. *Advances in Neural Information Processing Systems*, 34:21834–21846, 2021.
- [51] J. Zhu, Y. Yan, L. Zhao, M. Heimann, L. Akoglu, and D. Koutra. Beyond homophily in graph neural networks: Current limitations and effective designs. *Advances in Neural Information Processing Systems*, 33:7793–7804, 2020.
- [52] O. Shchur, M. Mumme, A. Bojchevski, and S. Günnemann. Pitfalls of graph neural network evaluation. In *NIPS workshop*, 2018.
- [53] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*. 2019.
- [54] M. Fey and J. E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [55] L. Biewald. Experiment tracking with weights and biases, 2020. Software available from wandb.com.

## Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section ??.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

### 1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? **[Yes]**
- (b) Did you describe the limitations of your work? **[Yes]**, in Section 6.
- (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** in the **Societal impact** paragraph in Section 6.
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**

### 2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
- (b) Did you include complete proofs of all theoretical results? **[Yes]** in Appendix A, Appendix B and Appendix C.

### 3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** Code and README in SM, dataloaders in code
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** Splits and hyperparameters provided in code zip
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[Yes]** Standard deviations are stated in results table
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** in appendix D

### 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? **[Yes]** datasets and standard libraries cited in appendix D
- (b) Did you mention the license of the assets? **[Yes]** industry standard libraries and benchmark datasets were used in accordance with licences
- (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** code provided in SM zip
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **[N/A]**
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[Yes]** no personal data is contained within benchmarking datasets

### 5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**

- 572 (b) Did you describe any potential participant risks, with links to Institutional Review  
573 Board (IRB) approvals, if applicable? [N/A]
- 574 (c) Did you include the estimated hourly wage paid to participants and the total amount  
575 spent on participant compensation? [N/A]

## A Proofs and additional details of Section 2

### A.1 Discussion on continuous Dirichlet energy and harmonic maps

In this subsection we briefly expand on the formulation of continuous Dirichlet energy in Section 2 to provide more context. Consider a smooth map  $f : (M, g) \rightarrow (N, h)$ , where  $N$  is usually a larger manifold we embed  $M$  into, and  $g, h$  are Riemannian metrics on domain and codomain respectively. The *Dirichlet energy* of  $f$  is defined by

$$\mathcal{E}(f, g, h) := \frac{1}{2} \int_M |df|_g^2 d\mu(g),$$

with  $|df|_g$  the norm of the Jacobian of  $f$  measured with respect to  $g$  and  $h$ . If  $(M, g)$  is standard Euclidean space  $\mathbb{R}^n$ ,  $N = \mathbb{R}^d$  and  $h$  is a constant positive semi-definite matrix, then we can rewrite the Dirichlet energy in a more familiar form as

$$\mathcal{E}(f, h) = \frac{1}{2} \int_{\mathbb{R}^n} \text{trace}(Df^\top h Df) d\mu = \frac{1}{2} \sum_{q,r=1}^d \sum_{j=1}^n \int_{\mathbb{R}^n} h_{qr} \partial_j f^q \partial_j f^r(x) dx.$$

The Dirichlet energy measures the smoothness of the map  $f$ , and indeed if  $h$  is the identity in  $\mathbb{R}^d$ , then we recover the classical definition

$$\mathcal{E}(f) = \frac{1}{2} \sum_{r=1}^d \int_{\mathbb{R}^n} \|\nabla f^r\|^2(x) dx.$$

**Gradient flow of Dirichlet energy.** Minimizers of  $\mathcal{E}$  - referred to as *harmonic maps* - are important objects in geometry: to mention a few, geodesics, minimal isometric immersions and maps  $f : M \rightarrow \mathbb{R}^d$  solving  $\Delta_g f = 0$  are all instances of harmonic maps. To identify such critical points, one computes the first variation of the energy  $\mathcal{E}$  along an arbitrary direction  $\partial_t f$ , which can be written as

$$d\mathcal{E}_f(\partial_t f) = - \int_M \langle \tau_g(f), \partial_t f \rangle_h d\mu(g).$$

for some tensor field  $\tau$  with explicit form

$$(\tau_{g_M}(f))^\alpha := \Delta_{g_M} f^\alpha + h_N \Gamma_{\beta\gamma}^\alpha \partial_i f^\beta \partial_j f^\gamma g_M^{ij},$$

for  $1 \leq \alpha \leq \dim(N)$ , with  $\{y^\alpha\}$  local coordinates on  $N$  and  $\Gamma_{\beta\gamma}^\alpha$  Christoffel symbols. It follows that harmonic maps are identified by the condition  $\tau_g(f) = 0$ . In [16], the pivotal idea of harmonic map flow – which has shaped much of modern research in geometric analysis – was introduced for the first time: in order to identify minimizers of  $\mathcal{E}$ , an input map  $f_0$  is evolved along the direction of (minus) the gradient of the energy  $\mathcal{E}$  leading to the dynamics

$$\partial_t f = \tau_g(f). \tag{14}$$

As a special case, when the target space is the classical Euclidean space one recovers the *heat equation* induced by the input Riemannian structure. We also note that when  $(M, g)$  is a surface representing an image and  $f : (u_1, u_2) \mapsto (u_1, u_2, \phi(u_1, u_2))$  with  $\phi$  a color map, then eq. (14) becomes

$$\partial_t \phi = \text{div}(C_g \nabla \phi), \tag{15}$$

with  $C_g$  a constant depending on the metric on  $M$ . If we now let  $g$  to depend on  $\phi$ , one can recover the celebrated Perona-Malik flow [26].

### A.2 Review of Kronecker product and properties of Laplacian kernel

**Additional notations and conventions used throughout the appendix.** Any graph  $G$  is taken to be *connected*. We order the eigenvalues of the graph Laplacian as  $0 = \lambda_0^\Delta \leq \lambda_1^\Delta \leq \dots \leq \lambda_{n-1}^\Delta = \rho_\Delta \leq 2$  with associated orthonormal basis of eigenvectors  $\{\phi_i^\Delta\}_{i=0}^{n-1}$  so that in particular we have  $\Delta \phi_0^\Delta = 0$ . Moreover, given a symmetric matrix  $\mathbf{B}$ , we generally denote the spectrum of  $\mathbf{B}$  by  $\text{spec}(\mathbf{B})$ . Finally, if we write  $\mathbf{F}(t)/\|\mathbf{F}(t)\|$  we always take the norm to be the Frobenius one and tacitly assume that the dynamics is s.t. the solution is not zero.

**Kronecker product.** In this subsection we summarize a few relevant notions pertaining the Kronecker product of matrices that are going to be applied throughout our spectral analysis of gradient flow equations for GNNs in both the continuous and discrete time setting.

Given a matricial equation of the form

$$\mathbf{Y} = \mathbf{A}\mathbf{X}\mathbf{B},$$

we can vectorize  $\mathbf{X}$  and  $\mathbf{Y}$  by stacking columns into  $\text{vec}(\mathbf{X})$  and  $\text{vec}(\mathbf{Y})$  respectively, and rewrite the previous system as

$$\text{vec}(\mathbf{Y}) = (\mathbf{B}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{X}). \quad (16)$$

If  $\mathbf{A}$  and  $\mathbf{B}$  are symmetric with spectra  $\text{spec}(\mathbf{A})$  and  $\text{spec}(\mathbf{B})$  respectively, then the spectrum of  $\mathbf{B} \otimes \mathbf{A}$  is given by  $\text{spec}(\mathbf{A}) \cdot \text{spec}(\mathbf{B})$ . Namely, if  $\mathbf{A}\mathbf{x} = \lambda^{\mathbf{A}}\mathbf{x}$  and  $\mathbf{B}\mathbf{y} = \lambda^{\mathbf{B}}\mathbf{y}$ , for  $\mathbf{x}$  and  $\mathbf{y}$  non-zero vectors, then  $\lambda^{\mathbf{B}}\lambda^{\mathbf{A}}$  is an eigenvalue of  $\mathbf{B} \otimes \mathbf{A}$  with eigenvector  $\mathbf{y} \otimes \mathbf{x}$ :

$$(\mathbf{B} \otimes \mathbf{A}) \mathbf{y} \otimes \mathbf{x} = (\lambda^{\mathbf{B}}\lambda^{\mathbf{A}}) \mathbf{y} \otimes \mathbf{x}. \quad (17)$$

One can also define the *Kronecker sum* of matrices  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{B} \in \mathbb{R}^{d \times d}$  as

$$\mathbf{A} \oplus \mathbf{B} := \mathbf{A} \otimes \mathbf{I}_d + \mathbf{I}_n \otimes \mathbf{B}, \quad (18)$$

with spectrum  $\text{spec}(\mathbf{A} \oplus \mathbf{B}) = \{\lambda^{\mathbf{A}} + \lambda^{\mathbf{B}} : \lambda^{\mathbf{A}} \in \text{spec}(\mathbf{A}), \lambda^{\mathbf{B}} \in \text{spec}(\mathbf{B})\}$ .

**Additional details on  $\mathcal{E}^{\text{Dir}}$  and the choice of Laplacian.** We recall that the classical graph Dirichlet energy  $\mathcal{E}^{\text{Dir}}$  is defined by

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}) = \frac{1}{2} \text{trace}(\mathbf{F}^\top \Delta \mathbf{F}),$$

where the (unusual) extra factor of  $\frac{1}{2}$  is to avoid rescaling the gradient flow by 2 – which is the more common convention. We can use the Kronecker product to rewrite the Dirichlet energy as

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}) = \frac{1}{2} \text{vec}(\mathbf{F})^\top (\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F}), \quad (19)$$

from which we immediately derive that  $\nabla_{\text{vec}(\mathbf{F})} \mathcal{E}^{\text{Dir}}(\mathbf{F}) = (\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F})$  – since  $\Delta$  is *symmetric* – and hence recover the gradient flow in eq. (3) leading to the graph heat equation across each channel.

Before we further comment on the characterizations of LFD and HFD dynamics, we review the main choices of graph Laplacian and the associated harmonic signals (i.e. how we can characterize the kernel spaces of the given Laplacian operator). Recall that throughout the appendix we always assume that the underlying graph  $G$  is *connected*. The symmetrically normalized Laplacian  $\Delta = \mathbf{I} - \mathbf{A}$  is symmetric, positive semi-definite with harmonic space of the form [14]

$$\ker(\Delta) := \text{span}(\mathbf{D}^{\frac{1}{2}} \mathbf{1}_n : \mathbf{1}_n = (1, \dots, 1)^\top). \quad (20)$$

This confirms that if a given GNN evolution  $\dot{\mathbf{F}}(t) = \text{GNN}_\theta(\mathbf{F}(t), t)$  with initial condition  $\mathbf{F}(0)$  over-smooths as per Definition 2.1 – i.e.  $\Delta \mathbf{f}^r(t) \rightarrow \mathbf{0}$  for  $t \rightarrow \infty$  for each column  $1 \leq r \leq d$  – then the only information persisting in the asymptotic regime is the degree and any dependence on the input features is lost, as studied in [30, 8]. A slightly different behaviour occurs if instead of  $\Delta$ , we consider the unnormalized Laplacian  $\mathbf{L} = \mathbf{D} - \mathbf{A}$  with kernel  $\text{span}(\mathbf{1}_n)$ , meaning that if  $\mathbf{L}\mathbf{f}^r(t) \rightarrow \mathbf{0}$  as  $t \rightarrow \infty$  for each  $1 \leq r \leq d$ , then any node would be embedded to a single point, hence making any separation task impossible. The same consequence applies to the random walk Laplacian  $\Delta_{\text{RW}} = \mathbf{I} - \mathbf{D}^{-1}\mathbf{A}$ . In particular, we note that generally a row-stochastic matrix is not symmetric – if it was, then this would in fact be doubly-stochastic – and the same applies to the random-walk Laplacian (a special exception is given by the class of *regular* graphs). In fact, in general any dynamical system governed by  $\Delta_{\text{RW}}$  (or simply  $\mathbf{D}^{-1}\mathbf{A}$ ) is not the gradient flow of an energy due to the lack of symmetry, as further confirmed below in eq. (22).

### A.3 Additional details on LFD and HFD characterizations

In this subsection we provide further details and justifications for Definition 2.1 and Definition 2.2. We first prove the following simple properties.

**Lemma A.1.** Assume we have a (continuous) process  $t \mapsto \mathbf{F}(t) \in \mathbb{R}^{n \times d}$ , for  $t \geq 0$ . The following equivalent characterizations hold:



- 648 (i)  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow 0$  for  $t \rightarrow \infty$  if and only if  $\Delta \mathbf{f}^r(t) \rightarrow \mathbf{0}$ , for  $1 \leq r \leq d$ .  
 649 (ii)  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)/\|\mathbf{F}(t)\|) \rightarrow \rho_\Delta/2$  for  $t \rightarrow \infty$  if and only if for any sequence  $t_j \rightarrow \infty$  there  
 650 exist a subsequence  $t_{j_k} \rightarrow \infty$  and a unit limit  $\mathbf{F}_\infty$  – depending on the subsequence – such  
 651 that  $\Delta \mathbf{f}_\infty^r = \rho_\Delta \mathbf{f}_\infty^r$ , for  $1 \leq r \leq d$ .

652 *Proof.* (i) Given  $\mathbf{F}(t) \in \mathbb{R}^{n \times d}$ , we can vectorize it and decompose it in the orthonormal basis  
 653  $\{\mathbf{e}_r \otimes \phi_i^\Delta : 1 \leq r \leq d, 0 \leq i \leq n-1\}$ , with  $\{\mathbf{e}_r\}_{r=1}^d$  canonical basis in  $\mathbb{R}^d$ , and write

$$\text{vec}(\mathbf{F}(t)) = \sum_{r,i} c_{r,i}(t) \mathbf{e}_r \otimes \phi_i^\Delta, \quad c_{r,i}(t) := \langle \text{vec}(\mathbf{F}(t)), \mathbf{e}_r \otimes \phi_i^\Delta \rangle.$$

654 We can then use eq. (19) to compute the Dirichlet energy as

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) = \frac{1}{2} \sum_{r=1}^d \sum_{i=0}^{n-1} c_{r,i}^2(t) \lambda_i^\Delta \equiv \frac{1}{2} \sum_{r=1}^d \sum_{i=1}^{n-1} c_{r,i}^2(t) \lambda_i^\Delta \geq \frac{1}{2} \text{gap}(\Delta) \sum_{r=1}^d \sum_{i=1}^{n-1} c_{r,i}^2(t),$$

655 where we have used the convention above that the eigenvector  $\phi_0^\Delta$  is in the kernel of  $\Delta$ . Therefore

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow 0 \iff \sum_{r=1}^d \sum_{i=1}^{n-1} c_{r,i}^2(t) \rightarrow 0, \quad t \rightarrow \infty,$$

656 which occurs if and only if

$$(\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F}(t)) = \sum_{r=1}^d \sum_{i=1}^{n-1} c_{r,i}(t) \lambda_i^\Delta \mathbf{e}_r \otimes \phi_i^\Delta \rightarrow 0.$$

657 (ii) The argument here is similar. Indeed we can write  $\mathbf{Q}(t) = \mathbf{F}(t)/\|\mathbf{F}(t)\|$  with  $\mathbf{Q}(t)$  a unit-norm  
 658 signal. Namely, we can vectorize and write

$$\text{vec}(\mathbf{Q}(t)) = \sum_{r,i} q_{r,i}(t) \mathbf{e}_r \otimes \phi_i^\Delta, \quad \sum_{r,i} q_{r,i}^2(t) = 1.$$

659 Then  $\mathcal{E}^{\text{Dir}}(\mathbf{Q}(t)) \rightarrow \rho_\Delta/2$  if and only if

$$\sum_{r,i} q_{r,i}^2(t) \lambda_i^\Delta \rightarrow \rho_\Delta, \quad t \rightarrow \infty,$$

660 which holds if and only if

$$\begin{aligned} \sum_r q_{r,\rho_\Delta}^2(t) &\rightarrow 1 \\ q_{r,i}^2(t) &\rightarrow 0, \quad i : \lambda_i^\Delta < \rho_\Delta, \end{aligned} \tag{21}$$

661 given the unit norm constraint. This is equivalent to the Rayleigh quotient of  $\mathbf{I}_d \otimes \Delta$  converging to its  
 662 maximal value  $\rho_\Delta$ . When this occurs, for any sequence  $t_j \rightarrow \infty$  we have that  $q_{r,i}^2(t_j) \leq 1$ , meaning  
 663 that we can extract a converging subsequence that due to eq. (21) will converge to a unit eigenvector  
 664  $\mathbf{Q}_\infty$  of  $\mathbf{I}_d \otimes \Delta$  satisfying  $(\mathbf{I}_d \otimes \Delta) \mathbf{Q}_\infty = \rho_\Delta \mathbf{Q}_\infty$ . Conversely assume for a contradiction that  
 665 there exists a sequence  $t_j \rightarrow \infty$  such that  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t_j)/\|\mathbf{F}(t_j)\|) < \rho_\Delta/2 - \epsilon$ , for some  $\epsilon > 0$ . Then  
 666 eq. (21) fails to be satisfied along the sequence, meaning that no subsequence converges to a unit  
 667 norm eigenvector  $\mathbf{F}_\infty$  of  $\mathbf{I}_d \otimes \Delta$  with associated eigenvalue  $\rho_\Delta$  which is a contradiction to our  
 668 assumption.

669 □

670 Before we address the formulation of low(high)-frequency-dominated dynamics, we solve explicitly  
 671 the system  $\dot{\mathbf{F}}(t) = \bar{\mathbf{A}} \mathbf{F}(t)$  in  $\mathbb{R}^{n \times d}$ , with some initial condition  $\mathbf{F}(0)$ . We can vectorize the equation  
 672 and solve  $\text{vec}(\mathbf{F}(t)) = (\mathbf{I}_d \otimes \bar{\mathbf{A}}) \text{vec}(\mathbf{F}(0))$ , meaning that

$$\text{vec}(\mathbf{F}(t)) = \sum_{r=1}^d \sum_{i=0}^{n-1} e^{(1-\lambda_i^\Delta)t} c_{r,i}(0) \mathbf{e}_r \otimes \phi_i^\Delta, \quad c_{r,i}(0) := \langle \text{vec}(\mathbf{F}(0)), \mathbf{e}_r \otimes \phi_i^\Delta \rangle.$$

673 Consider any initial condition  $\mathbf{F}(0)$  such that

$$\sum_{r=1}^d |c_{r,0}| = \sum_{r=1}^d \left| \langle \text{vec}(\mathbf{F}(0)), \mathbf{e}_r \otimes \phi_0^\Delta \rangle \right| > 0,$$

674 which is satisfied for each  $\text{vec}(\mathbf{F}(0)) \in \mathbb{R}^{nd} \setminus \mathcal{U}^\perp$ , where  $\mathcal{U}^\perp$  is the orthogonal complement of  
 675  $\mathbb{R}^d \otimes \text{span}(\phi_0^\Delta)$ . Since  $\mathcal{U}^\perp$  is a lower-dimensional subspace, its complement is dense. Accordingly  
 676 for a.e.  $\mathbf{F}(0)$ , we find that the solution satisfies

$$\|\text{vec}(\mathbf{F}(t))\|^2 = e^{2t} \left( \sum_{r=1}^d c_{r,0}^2 + \mathcal{O}(e^{-2\text{gap}(\Delta)t}) \right) = e^{2t} \left( \|P_{\ker(\Delta)}^\perp \text{vec}(\mathbf{F}(0))\|^2 + \mathcal{O}(e^{-2\text{gap}(\Delta)t}) \right),$$

677 with  $P_{\ker(\Delta)}^\perp$  the projection onto  $\mathbb{R}^d \otimes \ker(\Delta)$ . We see that the norm of the solution increases  
 678 exponentially, however the dominant term is given by the projection onto the lowest frequency signal  
 679 and in fact

$$\frac{\text{vec}(\mathbf{F}(t))}{\|\text{vec}(\mathbf{F}(t))\|} = \frac{P_{\ker(\Delta)}^\perp \text{vec}(\mathbf{F}(0)) + \mathcal{O}(e^{-\text{gap}(\Delta)t})(\mathbf{I} - P_{\ker(\Delta)}^\perp) \text{vec}(\mathbf{F}(0))}{\left( \|P_{\ker(\Delta)}^\perp \text{vec}(\mathbf{F}(0))\|^2 + \mathcal{O}(e^{-2\text{gap}(\Delta)t}) \right)^{\frac{1}{2}}} \rightarrow \text{vec}(\mathbf{F}_\infty),$$

680 such that  $(\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F}_\infty) = \mathbf{0}$  which means  $\Delta \mathbf{f}_\infty^r = \mathbf{0}$ , for each column  $1 \leq r \leq d$ . Equivalently,  
 681 one can compute  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)/\|\mathbf{F}(t)\|)$  and conclude that the latter quantity converges to zero as  
 682  $t \rightarrow \infty$  by the very same argument.

683 In fact, this motivates further the nomenclature LFD and HFD. Without loss of generality we  
 684 focus now on the high-frequency case. Assume that we have a HFD dynamics  $t \mapsto \mathbf{F}(t)$ ,  
 685 i.e.  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)/\|\mathbf{F}(t)\|) \rightarrow \rho_\Delta/2$ , then we can vectorize the solution and write  $\text{vec}(\mathbf{F}(t)) =$   
 686  $\|\mathbf{F}(t)\| \text{vec}(\mathbf{Q}(t))$ , for some time-dependent unit vector  $\text{vec}(\mathbf{Q}(t)) \in \mathbb{R}^{nd}$ :

$$\text{vec}(\mathbf{Q}(t)) = \sum_{r,i} q_{r,i}(t) \mathbf{e}_r \otimes \phi_i^\Delta, \quad \sum_{r,i} q_{r,i}^2(t) = 1.$$

687 By Lemma A.1 and more explicitly eq. (21), we derive that the coefficients  $\{q_{r,\rho_\Delta}\}$  associated with  
 688 the eigenvectors  $\mathbf{e}_r \otimes \phi_{\rho_\Delta}^\Delta$  are dominant in the evolution hence justifying the name *high-frequency*  
 689 *dominated* dynamics.

690 We note that the next result covers Lemma 2.3.

691 **Lemma A.2.** Consider a dynamical system  $\dot{\mathbf{F}}(t) = \text{GNN}_\theta(\mathbf{F}(t), t)$ , with initial condition  $\mathbf{F}(0)$ .

692 (i)  $\text{GNN}_\theta$  is LFD if and only if  $(\mathbf{I}_d \otimes \Delta) \frac{\text{vec}(\mathbf{F}(t))}{\|\mathbf{F}(t)\|} \rightarrow \mathbf{0}$  if and only if for each sequence  
 693  $t_j \rightarrow \infty$  there exist a subsequence  $t_{j_k} \rightarrow \infty$  and  $\mathbf{F}_\infty$  (depending on the subsequence) s.t.  
 694  $\frac{\mathbf{F}(t_{j_k})}{\|\mathbf{F}(t_{j_k})\|} \rightarrow \mathbf{F}_\infty$  satisfying  $\Delta \mathbf{f}_\infty^r = \mathbf{0}$ , for each  $1 \leq r \leq d$ .

695 (ii)  $\text{GNN}_\theta$  is HFD if and only if for each sequence  $t_j \rightarrow \infty$  there exist a subsequence  $t_{j_k} \rightarrow \infty$   
 696 and  $\mathbf{F}_\infty$  (depending on the subsequence) s.t.  $\frac{\mathbf{F}(t_{j_k})}{\|\mathbf{F}(t_{j_k})\|} \rightarrow \mathbf{F}_\infty$  satisfying  $\Delta \mathbf{f}_\infty^r = \rho_\Delta \mathbf{f}_\infty^r$ ,  
 697 for each  $1 \leq r \leq d$ .

698 *Proof.* (i) Since  $\Delta \mathbf{f}^r(t) \rightarrow \mathbf{0}$  for each  $1 \leq r \leq d$  if and only if  $(\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F}(t)) \rightarrow \mathbf{0}$ , we  
 699 conclude that the dynamics is LFD if and only if  $(\mathbf{I}_d \otimes \Delta) \frac{\text{vec}(\mathbf{F}(t))}{\|\mathbf{F}(t)\|} \rightarrow \mathbf{0}$  due to (i) in Lemma A.1.

700 Consider a sequence  $t_j \rightarrow \infty$ . Since  $\text{vec}(\mathbf{F}(t_j))/\|\mathbf{F}(t_j)\|$  is a bounded sequence we can extract  
 701 a converging subsequence  $t_{j_k}$ :  $\text{vec}(\mathbf{F}(t_{j_k}))/\|\mathbf{F}(t_{j_k})\| \rightarrow \text{vec}(\mathbf{F}_\infty)$ . If the dynamics is LFD, then  
 702  $(\mathbf{I}_d \otimes \Delta) \frac{\text{vec}(\mathbf{F}(t_{j_k}))}{\|\mathbf{F}(t_{j_k})\|} \rightarrow \mathbf{0}$  and hence we conclude that  $\text{vec}(\mathbf{F}_\infty) \in \ker(\mathbf{I}_d \otimes \Delta)$ . Conversely, assume

703 that for any sequence  $t_j \rightarrow \infty$  there exists a subsequence  $t_{j_k}$  and  $\mathbf{F}_\infty$  such that  $\frac{\mathbf{F}(t_{j_k})}{\|\mathbf{F}(t_{j_k})\|} \rightarrow \mathbf{F}_\infty$   
 704 satisfying  $\Delta \mathbf{f}_\infty^r = \mathbf{0}$ , for each  $1 \leq r \leq d$ . If for a contradiction we had  $\varepsilon > 0$  and  $t_j \rightarrow \infty$  such that  
 705  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t_j)/\|\mathbf{F}(t_j)\|) \geq \varepsilon$  – for  $j$  large enough – then by (i) in Lemma A.1 there exist  $1 \leq r \leq d$ ,  
 706  $i > 0$  and a subsequence  $t_{j_k}$  satisfying

$$\left| \left\langle \frac{\text{vec}(\mathbf{F}(t_{j_k}))}{\|\mathbf{F}(t_{j_k})\|}, \mathbf{e}_r \otimes \phi_i^\Delta \right\rangle \right| > \delta(\varepsilon) > 0,$$

707 meaning that there is no subsequence of  $\{t_{j_k}\}$  s.t.  $(\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F}(t_{j_k})) / \|\mathbf{F}(t_{j_k})\| \rightarrow \mathbf{0}$ , providing  
 708 a contradiction.

709 (ii) This is equivalent to (ii) in Lemma A.1.

710 □

711 **Remark.** We note that in Lemma 2.3 an LFD dynamics does not necessarily mean that the normalized  
 712 solution converges to the kernel of  $\mathbf{I}_d \otimes \Delta$  – i.e. one in general has always to pass to subsequences.  
 713 Indeed, we can consider the simple example  $t \mapsto \text{vec}(\mathbf{F}(t)) := \cos(t) \mathbf{e}_{\bar{r}} \otimes \phi_0^\Delta$ , for some  $1 \leq \bar{r} \leq d$ ,  
 714 which satisfies  $\Delta \mathbf{f}^r(t) = \mathbf{0}$  for each  $r$ , but it is not a convergent function due to its oscillatory nature.  
 715 Same argument applies to HFD.

#### 716 A.4 Details and proofs on $\mathcal{E}_{\mathbf{W}}^{\text{Dir}}$ and its gradient flow

717 By direct computation one verifies that the definition in eq. (4) can be equivalently written as

$$\mathcal{E}_{\mathbf{W}}^{\text{Dir}}(\mathbf{F}) = \frac{1}{2} \langle \text{vec}(\mathbf{F}), (\mathbf{W}^\top \mathbf{W} \otimes \Delta) \text{vec}(\mathbf{F}) \rangle,$$

718 from which we immediately derive  $\nabla_{\text{vec}(\mathbf{F})} \mathcal{E}_{\mathbf{W}}^{\text{Dir}}(\text{vec}(\mathbf{F})) = (\mathbf{W}^\top \mathbf{W} \otimes \Delta) \text{vec}(\mathbf{F})$  which proves  
 719 eq. (5). We can now address the proof of Proposition 2.4.

720 *Proof of Proposition 2.4.* We can vectorize the gradient flow system in eq. (5) and use the spectral  
 721 characterization of  $\mathbf{W}^\top \mathbf{W} \otimes \Delta$  in eq. (17) to write the solution explicitly as

$$\text{vec}(\mathbf{F}(t)) = \sum_{r,i} e^{-(\lambda_r^{\mathbf{W}} \lambda_i^\Delta) t} c_{r,i}(0) \phi_r^{\mathbf{W}} \otimes \phi_i^\Delta,$$

722 where  $\{\lambda_r^{\mathbf{W}}\}_r = \text{spec}(\mathbf{W}^\top \mathbf{W}) \subset \mathbb{R}_{\geq 0}$  with associated basis of orthonormal eigenvectors given by  
 723  $\{\phi_r^{\mathbf{W}}\}_r$ . Then

$$\begin{aligned} \mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) &= \frac{1}{2} \langle \text{vec}(\mathbf{F}(t)), (\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F}(t)) \rangle = \frac{1}{2} \sum_{r,i} e^{-2t(\lambda_r^{\mathbf{W}} \lambda_i^\Delta)} c_{r,i}^2(0) \lambda_i^\Delta \\ &= \frac{1}{2} \sum_{r: \lambda_r^{\mathbf{W}}=0, i} c_{r,i}^2(0) \lambda_i^\Delta + \frac{1}{2} \sum_{r: \lambda_r^{\mathbf{W}}>0, i>0} c_{r,i}^2(0) e^{-2t(\lambda_r^{\mathbf{W}} \lambda_i^\Delta)} \lambda_i^\Delta \\ &= \mathcal{E}^{\text{Dir}}((P_{\mathbf{W}}^{\text{ker}} \otimes \mathbf{I}_n) \text{vec}(\mathbf{F}(0))) + \frac{1}{2} \sum_{r: \lambda_r^{\mathbf{W}}>0, i>0} c_{r,i}^2(0) e^{-2t(\lambda_r^{\mathbf{W}} \lambda_i^\Delta)} \lambda_i^\Delta \\ &\leq \mathcal{E}^{\text{Dir}}((P_{\mathbf{W}}^{\text{ker}} \otimes \mathbf{I}_n) \text{vec}(\mathbf{F}(0))) + \frac{\rho \Delta}{2} e^{-2t \text{gap}(\mathbf{W}^\top \mathbf{W}) \text{gap}(\Delta)} \|\mathbf{F}(0)\|^2, \end{aligned}$$

724 where we recall that  $P_{\mathbf{W}}^{\text{ker}}$  is the projection onto  $\ker(\mathbf{W}^\top \mathbf{W})$  and that by convention the index  $i = 0$   
 725 is associated with the lowest graph frequency  $\lambda_0^\Delta = 0$  – by assumption G is connected. This proves  
 726 that the dynamics is in fact smoothing as per Definition 2.1. By the very same argument we find that

$$\text{vec}(\mathbf{F}(t)) \rightarrow (\mathbf{I}_d \otimes P_{\Delta}^{\text{ker}}) \text{vec}(\mathbf{F}(0)) + (P_{\mathbf{W}}^{\text{ker}} \otimes \mathbf{I}_n) \text{vec}(\mathbf{F}(0)), \quad t \rightarrow \infty,$$

727 with  $P_{\Delta}^{\text{ker}}$  the orthogonal projection onto  $\ker \Delta$  – the other terms decay exponentially to zero. We  
 728 first focus on the first quantity, which we can write as

$$(\mathbf{I}_d \otimes P_{\Delta}^{\text{ker}}) \text{vec}(\mathbf{F}(0)) = \sum_r c_{r,0}(0) \phi_r^{\mathbf{W}} \otimes \phi_0^\Delta,$$

729 which has matrix representation  $\phi_0^\Delta \phi_\infty^\top \in \mathbb{R}^{n \times d}$  with

$$\phi_\infty := \sum_r c_{r,0}(0) \phi_r^{\mathbf{W}}.$$

730 By eq. (20) we deduce that the  $i$ -th row of  $\phi_0^\Delta \phi_\infty^\top \in \mathbb{R}^{n \times d}$  is the  $d$ -dimensional vector  $\sqrt{d_i} \phi_\infty$ . We  
 731 now focus on the term

$$(P_{\mathbf{W}}^{\text{ker}} \otimes \mathbf{I}_n) \text{vec}(\mathbf{F}(0)) = \sum_{r: \lambda_r^{\mathbf{W}}=0, j} c_{r,j}(0) \phi_r^{\mathbf{W}} \otimes \phi_j^\Delta$$

732 which has matrix representation  $\sum_{r:\lambda_r^{\mathbf{W}}=0,j} c_{r,j}(0)\phi_j^{\Delta}(\phi_r^{\mathbf{W}})^{\top}$ . In particular, the  $i$ -th row is given by

$$\sum_{r:\lambda_r^{\mathbf{W}}=0,j} c_{r,j}(0)(\phi_j^{\Delta})_i \phi_r^{\mathbf{W}} = P_{\mathbf{W}}^{\ker} \mathbf{f}_i(0).$$

733 This completes the proof of Proposition 2.4.  $\square$

## 734 B Proofs and additional details of Section 3

### 735 B.1 Spectral analysis of the channel-mixing: the continuous case

736 Consider the generalized energy  $\mathcal{E}^{\text{tot}}$  in eq. (7). We can use vectorization to rewrite it as

$$\mathcal{E}^{\text{tot}}(\text{vec}(\mathbf{F})) = \frac{1}{2} \langle \text{vec}(\mathbf{F}), (\mathbf{\Omega} \otimes \mathbf{I}_n) \text{vec}(\mathbf{F}) \rangle - \frac{1}{2} \langle \text{vec}(\mathbf{F}), (\mathbf{W} \otimes \bar{\mathbf{A}}) \text{vec}(\mathbf{F}) \rangle,$$

737 from which the gradient flow in eq. (8) follows. In particular, given a system as in eq. (8):

$$\text{vec}(\dot{\mathbf{F}}(t)) = -(\mathbf{\Omega} \otimes \mathbf{I}_n) \text{vec}(\mathbf{F}(t)) + (\mathbf{W} \otimes \bar{\mathbf{A}}) \text{vec}(\mathbf{F}(t)),$$

738 if this is the gradient flow of  $\mathbf{F} \mapsto \mathcal{E}^{\text{tot}}(\mathbf{F})$ , then we would have

$$\nabla_{\text{vec}(\mathbf{F})}^2 \mathcal{E}^{\text{tot}}(\mathbf{F}) = \mathbf{\Omega} \otimes \mathbf{I}_n - \mathbf{W} \otimes \bar{\mathbf{A}}, \quad (22)$$

739 which must be symmetric due to the Hessian of a function being symmetric. The latter means

$$(\mathbf{\Omega}^{\top} - \mathbf{\Omega}) \otimes \mathbf{I}_n = (\mathbf{W}^{\top} - \mathbf{W}) \otimes \bar{\mathbf{A}},$$

740 which is satisfied if and only if both  $\mathbf{\Omega}$  and  $\mathbf{W}$  are *symmetric*. This shows that eq. (8) *is the gradient*  
741 *flow of  $\mathcal{E}^{\text{tot}}$  if and only if  $\mathbf{\Omega}$  and  $\mathbf{W}$  are symmetric.*

742 We now rely on the spectral decomposition of  $\mathbf{W}$  to rewrite  $\mathcal{E}^{\text{tot}}$  explicitly in terms of attractive  
743 and repulsive interactions. If we have a spectral decomposition  $\mathbf{W} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^{\top}$ , we can separate the  
744 positive eigenvalues from the negative ones and write

$$\mathbf{W} = \mathbf{U} \mathbf{\Lambda}_+ \mathbf{U}^{\top} + \mathbf{U} \mathbf{\Lambda}_- \mathbf{U}^{\top} := \mathbf{W}_+ - \mathbf{W}_-.$$

745 Since  $\mathbf{W}_+ \succeq 0, \mathbf{W}_- \succeq 0$ , we can use the Choleski decomposition to write  $\mathbf{W}_+ = \mathbf{\Theta}_+^{\top} \mathbf{\Theta}_+$  and  
746  $\mathbf{W}_- = \mathbf{\Theta}_-^{\top} \mathbf{\Theta}_-$  with  $\mathbf{\Theta}_+, \mathbf{\Theta}_- \in \mathbb{R}^{d \times d}$ . Equation (9) follows then by direct computation: namely

$$\begin{aligned} \mathcal{E}^{\text{tot}}(\mathbf{F}) &= \frac{1}{2} \sum_i \langle \mathbf{f}_i, \mathbf{\Omega} \mathbf{f}_i \rangle - \frac{1}{2} \sum_{i,j} \bar{a}_{ij} \langle \mathbf{f}_i, \mathbf{W} \mathbf{f}_j \rangle \\ &= \frac{1}{2} \sum_i \langle \mathbf{f}_i, (\mathbf{\Omega} - \mathbf{W}) \mathbf{f}_i \rangle + \frac{1}{2} \sum_i \langle \mathbf{f}_i, \mathbf{W} \mathbf{f}_i \rangle - \frac{1}{2} \sum_{i,j} \bar{a}_{ij} \langle \mathbf{\Theta}_+ \mathbf{f}_i, \mathbf{\Theta}_+ \mathbf{f}_j \rangle + \frac{1}{2} \sum_{i,j} \bar{a}_{ij} \langle \mathbf{\Theta}_- \mathbf{f}_i, \mathbf{\Theta}_- \mathbf{f}_j \rangle \\ &= \frac{1}{2} \sum_i \langle \mathbf{f}_i, (\mathbf{\Omega} - \mathbf{W}) \mathbf{f}_i \rangle + \frac{1}{4} \sum_{i,j} \|\mathbf{\Theta}_+ (\nabla \mathbf{F})_{ij}\|^2 - \frac{1}{4} \sum_{i,j} \|\mathbf{\Theta}_- (\nabla \mathbf{F})_{ij}\|^2, \end{aligned}$$

747 where we have used that  $\sum_{i,j} \frac{1}{d_i} \|\mathbf{\Theta}_+ \mathbf{f}_i\|^2 = \sum_i \|\mathbf{\Theta}_+ \mathbf{f}_i\|^2$ .

748 *Proof of Proposition 3.1.* Once we compute the spectrum of  $\mathbf{W} \otimes \bar{\mathbf{A}}$  via eq. (17), we can write the  
749 solution as – recall that  $\bar{\mathbf{A}} = \mathbf{I}_n - \mathbf{\Delta}$  so we can rephrase the eigenvalues of  $\bar{\mathbf{A}}$  in terms of the  
750 eigenvalues of  $\mathbf{\Delta}$ :

$$\text{vec}(\mathbf{F}(t)) = \sum_{r,i} e^{\lambda_r^{\mathbf{W}}(1-\lambda_i^{\mathbf{\Delta}})t} c_{r,i}(0) \phi_r^{\mathbf{W}} \otimes \phi_i^{\mathbf{\Delta}},$$

751 with  $\mathbf{W} \phi_r^{\mathbf{W}} = \lambda_r^{\mathbf{W}} \phi_r^{\mathbf{W}}$ , for  $1 \leq r \leq d$ , where  $\{\phi_r^{\mathbf{W}}\}_r$  is an orthonormal basis of eigenvectors in  
752  $\mathbb{R}^d$ . We can then calculate the Dirichlet energy along the solution as

$$\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) = \frac{1}{2} \langle \text{vec}(\mathbf{F}(t)), (\mathbf{I}_d \otimes \mathbf{\Delta}) \text{vec}(\mathbf{F}(t)) \rangle = \frac{1}{2} \sum_{r,i} e^{2\lambda_r^{\mathbf{W}}(1-\lambda_i^{\mathbf{\Delta}})t} c_{r,i}^2(0) \lambda_i^{\mathbf{\Delta}}.$$

753 We now consider two cases:

- If  $\lambda_r^{\mathbf{W}} > 0$ , then  $\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}) \leq \lambda_+^{\mathbf{W}}$ .
- If  $\lambda_r^{\mathbf{W}} < 0$ , then  $\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}) \leq |\lambda_-^{\mathbf{W}}|(\rho_{\Delta} - 1) := \rho_-$ , with eigenvectors  $\phi_r^{\mathbf{W}} \otimes \phi_{\rho_{\Delta}}^{\Delta}$  for each  $r$  s.t.  $\mathbf{W}\phi_r^{\mathbf{W}} = \lambda_r^{\mathbf{W}}\phi_r^{\mathbf{W}}$  – without loss of generality we can assume that  $\rho_{\Delta}$  is a simple eigenvalue for  $\Delta$ . In particular, if  $\lambda_r^{\mathbf{W}} < 0$  and  $\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}) < \rho_-$ , then

$$\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}) < \max\{|\lambda_-^{\mathbf{W}}|(\lambda_{n-2}^{\Delta} - 1), |\lambda_{-,2}^{\mathbf{W}}|(\rho_{\Delta} - 1)\},$$

where  $\lambda_{-,2}^{\mathbf{W}}$  is the second most negative eigenvalue of  $\mathbf{W}$  and  $\lambda_{n-2}^{\Delta}$  is the second largest eigenvalue of  $\Delta$ . In particular, we can write

$$\lambda_{n-2}^{\Delta} = \rho_{\Delta} - \text{gap}(\rho_{\Delta}\mathbf{I}_n - \Delta), \quad |\lambda_{-,2}^{\mathbf{W}}| = |\lambda_-^{\mathbf{W}}| - \text{gap}(|\lambda_-^{\mathbf{W}}|\mathbf{I}_d + \mathbf{W}). \quad (23)$$

From (i) and (ii) we derive that if  $\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}) \neq \rho_-$ , then

$$\begin{aligned} \lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}) - \rho_- &< -\min\{\rho_- - \lambda_+^{\mathbf{W}}, \rho_- - |\lambda_-^{\mathbf{W}}|(\lambda_{n-2}^{\Delta} - 1), \rho_- - |\lambda_{-,2}^{\mathbf{W}}|(\rho_{\Delta} - 1)\} \\ &= -\min\{\rho_- - \lambda_+^{\mathbf{W}}, |\lambda_-^{\mathbf{W}}|\text{gap}(\rho_{\Delta}\mathbf{I} - \Delta), \text{gap}(|\lambda_-^{\mathbf{W}}|\mathbf{I} + \mathbf{W})(\rho_{\Delta} - 1)\} = -\epsilon_{\text{HFD}}, \end{aligned} \quad (24)$$

where we have used eq. (23). Accordingly, if  $\rho_- > \lambda_+^{\mathbf{W}}$ , then

$$\begin{aligned} \mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) &= e^{2t\rho_-} \left( \frac{\rho_{\Delta}}{2} \sum_{r:\lambda_r^{\mathbf{W}}=\lambda_-^{\mathbf{W}}} c_{r,\rho_{\Delta}}^2(0) + \frac{1}{2} \sum_{r,i:\lambda_r^{\mathbf{W}}(1-\lambda_i^{\Delta})\neq\rho_-} e^{2(\lambda_r^{\mathbf{W}}(1-\lambda_i^{\Delta})-\rho_-)t} c_{r,i}^2(0) \right) \\ &= e^{2t\rho_-} \left( \frac{\rho_{\Delta}}{2} \|P_{\mathbf{W}}^{\rho_-} \mathbf{F}(0)\|^2 + \mathcal{O}(e^{-2t\epsilon_{\text{HFD}}}) \right). \end{aligned}$$

By the same argument we can factor out the dominant term and derive the following limit for  $t \rightarrow \infty$  and for a.e.  $\mathbf{F}(0)$  since  $P_{\mathbf{W}}^{\rho_-} \text{vec}(\mathbf{F}(0)) = \mathbf{0}$  only if  $\text{vec}(\mathbf{F}(0))$  belongs to a lower dimensional subspace of  $\mathbb{R}^{nd}$ :

$$\frac{\text{vec}(\mathbf{F}(t))}{\text{vec}(\mathbf{F}(t))} = \frac{P_{\mathbf{W}}^{\rho_-} \text{vec}(\mathbf{F}(0)) + \mathcal{O}(e^{-\epsilon_{\text{HFD}}t})((\mathbf{I} - P_{\mathbf{W}}^{\rho_-})\text{vec}(\mathbf{F}(0)))}{(\|P_{\mathbf{W}}^{\rho_-} \text{vec}(\mathbf{F}(0))\|^2 + \mathcal{O}(e^{-2\epsilon_{\text{HFD}}t}))^{\frac{1}{2}}} \rightarrow \frac{P_{\mathbf{W}}^{\rho_-} \text{vec}(\mathbf{F}(0))}{\|P_{\mathbf{W}}^{\rho_-} \text{vec}(\mathbf{F}(0))\|},$$

where the latter is a unit vector  $\text{vec}(\mathbf{F}_{\infty})$  satisfying  $(\mathbf{I}_d \otimes \Delta)\text{vec}(\mathbf{F}_{\infty}) = \rho_{\Delta}\text{vec}(\mathbf{F}_{\infty})$ , which completes the proof.  $\square$

## B.2 Propagating with $-\Delta$ : a perspective in terms of channel-mixing spectrum

In this subsection we briefly review the special case of eq. (8) where  $\Omega = \mathbf{W}$ , and comment on why we generally expect a framework where the propagation is governed by the graph vector field  $\mathbf{A}$  to be more flexible than one with  $-\Delta$ . If  $\Omega = \mathbf{W}$ , the gradient flow in eq. (8) becomes

$$\dot{\mathbf{F}}(t) = -\Delta\mathbf{F}(t)\mathbf{W}. \quad (25)$$

We note that once vectorized, the solution to the dynamical system can be written as

$$\text{vec}(\mathbf{F}(t)) = \sum_{r=1}^d \sum_{i=0}^{n-1} e^{-\lambda_r^{\mathbf{W}}\lambda_i^{\Delta}t} c_{r,i}(0) \phi_r^{\mathbf{W}} \otimes \phi_i^{\Delta}.$$

In particular, we immediately deduce the following counterpart to Proposition 3.1

**Corollary B.1.** *If  $\text{spec}(\mathbf{W}) \cap \mathbb{R}_- \neq \emptyset$ , then eq. (25) is HFD for a.e.  $\mathbf{F}(0)$ .*

Differently from eq. (8) the lowest frequency component is *always preserved independent of the spectrum of  $\mathbf{W}$* . This means that the system cannot learn eigenvalues of  $\mathbf{W}$  to either magnify or suppress the low-frequency projection. In contrast, this can be done if  $\Omega = \mathbf{0}$ , or equivalently one replaces  $-\Delta$  with  $\mathbf{A}$  providing a further justification in terms of the interaction between graph spectrum and channel-mixing spectrum for why graph convolutional models use the normalized adjacency rather than the Laplacian for propagating messages [27].

### 780 B.3 A more general family of energies: gradient flow with non-linear activations

781 Consider a more general pairwise energy including a non-linear differentiable activation map  $\sigma$  of the  
782 form

$$\mathcal{E}_{\sigma, \mathbf{W}}^{\text{pair}}(\mathbf{F}) = \frac{1}{2} \sum_{i,j} \bar{a}_{i,j} \sigma(\mathbf{f}_i, \mathbf{W} \mathbf{f}_j).$$

783 We temporarily assume that  $\Omega = \mathbf{0}$ . The gradient flow follows from direct computation:

$$\dot{\mathbf{F}}(t) = \mathcal{A}_{\sigma}(\mathbf{F}(t)) \mathbf{F}(t) \mathbf{W}, \quad (\mathcal{A}_{\sigma}(\mathbf{F}(t)))_{ij} := \bar{a}_{ij} \sigma'(\mathbf{f}_i, \mathbf{W} \mathbf{f}_j). \quad (26)$$

784 In particular, we see that the non-linear activations in general may induce a type of attention mech-  
785 anism where the diffusion along edges is controlled by the derivative of  $\sigma$  evaluated on the inner  
786 product of features induced by  $\mathbf{W}$ . A similar structure is investigated in [17]. We also observe  
787 that analogous conclusions can be deduced if  $\Omega \neq \mathbf{0}$  and the external energy term  $\mathcal{E}_{\Omega}^{\text{ext}}$  includes a  
788 non-linear activation map  $\sigma$  as in the pairwise contribution.

## 789 C Proofs and additional details of Section 4

790 We first explicitly report here the expansion of the discrete gradient flow in eq. (11) after  $m$  layers to  
791 further highlight how this is not equivalent to a single linear layer with a message passing matrix  $\bar{\mathbf{A}}^m$   
792 as for SGCN [43]. For simplicity we suppress the source term.

$$\begin{aligned} \mathbf{F}(t + \tau) &= \mathbf{F}(t) + \tau (-\mathbf{F}(t) \Omega + \bar{\mathbf{A}} \mathbf{F}(t) \mathbf{W}) \\ \text{vec}(\mathbf{F}(t + \tau)) &= (\mathbf{I}_{nd} + \tau (-\Omega \otimes \mathbf{I}_n + \mathbf{W} \otimes \bar{\mathbf{A}})) \text{vec}(\mathbf{F}(t)) \\ \text{vec}(\mathbf{F}(m\tau)) &= \sum_{k=0}^m \binom{m}{k} \tau^k (-\Omega \otimes \mathbf{I}_n + \mathbf{W} \otimes \bar{\mathbf{A}})^k \text{vec}(\mathbf{F}(0)) \end{aligned} \quad (27)$$

793 and we see how the message passing matrix  $\bar{\mathbf{A}}$  actually enters the expansion after  $m$  layers with each  
794 power  $0 \leq k \leq m$ . This is not surprising, given that *we are discretizing a linear dynamical system,*  
795 *meaning that we are approximating an exponential matrix.*

### 796 C.1 Comparison with continuous GNNs: details and proofs

797 We prove the following result which covers Proposition 4.1.

798 *Proof of Proposition 4.1.* We structure the proof by following the numeration in the statement.

799 (i) From direct computation we find

$$\begin{aligned} \frac{d\mathcal{E}^{\text{Dir}}(\mathbf{F}(t))}{dt} &= \frac{1}{2} \frac{d}{dt} (\langle \text{vec}(\mathbf{F}(t)), (\mathbf{I}_d \otimes \Delta) \text{vec}(\mathbf{F}(t)) \rangle) \\ &= -\langle \text{vec}(\mathbf{F}(t)), (\mathbf{K}^{\top}(t) \mathbf{K}(t) \otimes \Delta^2) \text{vec}(\mathbf{F}(t)) \rangle \leq 0, \end{aligned}$$

800 since  $\mathbf{K}^{\top}(t) \mathbf{K}(t) \otimes \Delta^2 \succeq 0$ . Note that we have used that  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD}$ .

801 (ii) We consider the dynamical system

$$\dot{\mathbf{F}}_{\text{CGNN}}(t) = -\Delta \mathbf{F}(t) + \mathbf{F}(t) \tilde{\Omega} + \mathbf{F}(0).$$

802 We can write  $\text{vec}(\mathbf{F}(t)) = \sum_{r,i} c_{r,i}(t) \phi_r^{\tilde{\Omega}} \otimes \phi_i^{\Delta}$ , leading to the system

$$\dot{c}_{r,i}(t) = (\lambda_r^{\tilde{\Omega}} - \lambda_i^{\Delta}) c_{r,i}(t) + c_{r,i}(0), \quad 0 \leq i \leq n-1, \quad 1 \leq r \leq d.$$

803 We can solve explicitly the system as

$$\begin{aligned} c_{r,i}(t) &= c_{r,i}(0) \left( e^{(\lambda_r^{\tilde{\Omega}} - \lambda_i^{\Delta})t} \left( 1 + \frac{1}{\lambda_r^{\tilde{\Omega}} - \lambda_i^{\Delta}} \right) - \frac{1}{\lambda_r^{\tilde{\Omega}} - \lambda_i^{\Delta}} \right), \quad \text{if } \lambda_r^{\tilde{\Omega}} \neq \lambda_i^{\Delta} \\ c_{r,i}(t) &= c_{r,i}(0)(1+t), \quad \text{otherwise.} \end{aligned}$$

804 We see now that for a.e.  $\mathbf{F}(0)$  the projection  $(\mathbf{I}_d \otimes \phi_{\rho_\Delta}^\Delta (\phi_{\rho_\Delta}^\Delta)^\top) \text{vec}(\mathbf{F}(t))$  is never the dominant  
805 term. In fact, if there exists  $r$  s.t.  $\lambda_r^{\tilde{\Omega}} \geq \rho_\Delta$ , then  $\lambda_r^{\tilde{\Omega}} - \lambda_i^\Delta > \lambda_r^{\tilde{\Omega}} - \rho_\Delta$ , for any other non-maximal  
806 graph Laplacian eigenvalue. It follows that there is *no*  $\tilde{\Omega}$  s.t. the normalized solution maximizes the  
807 Rayleigh quotient of  $\mathbf{I}_d \otimes \Delta$ , proving that CGNN is never HFD.

808 If we have no source, then the CGNN equation becomes

$$\dot{\mathbf{F}}(t) = -\Delta \mathbf{F}(t) + \mathbf{F}(t) \tilde{\Omega} \iff \text{vec}(\dot{\mathbf{F}}(t)) = (\tilde{\Omega} \oplus (-\Delta)) \text{vec}(\mathbf{F}(t)),$$

809 using the Kronecker sum notation in eq. (18). It follows that we can write the vectorized solution in  
810 the basis  $\{\phi_r^{\tilde{\Omega}} \otimes \phi_i^\Delta\}_{r,i}$  as

$$\begin{aligned} \text{vec}(\mathbf{F}(t)) = & e^{\lambda_+^{\tilde{\Omega}} t} \left( \sum_{r: \lambda_r^{\tilde{\Omega}} = \lambda_+^{\tilde{\Omega}}} c_{r,0}(0) \phi_r^{\tilde{\Omega}} \otimes \phi_0^\Delta + \mathcal{O}(e^{-\text{gap}(\lambda_+^{\tilde{\Omega}} - \tilde{\Omega})t}) \sum_{r: \lambda_r^{\tilde{\Omega}} < \lambda_+^{\tilde{\Omega}}} c_{r,0}(0) \phi_r^{\tilde{\Omega}} \otimes \phi_0^\Delta \right) \\ & + e^{\lambda_+^{\tilde{\Omega}} t} \left( \mathcal{O}(e^{-\text{gap}(\Delta)t}) \left( \sum_{r,i>0} c_{r,i}(0) \phi_r^{\tilde{\Omega}} \otimes \phi_i^\Delta \right) \right), \end{aligned}$$

811 meaning that the dominant term is given by the lowest frequency component and in fact, if we  
812 normalize we find  $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)/\|\mathbf{F}(t)\|) \leq e^{-\text{gap}(\Delta)t}$ .

813 (iii) Finally we consider the dynamical system induced by linear GRAND

$$\dot{\mathbf{F}}_{\text{GRAND}}(t) = -\Delta_{\text{RW}} \mathbf{F}(t) = -(\mathbf{I} - \mathcal{A}(\mathbf{F}(0))) \mathbf{F}(t).$$

814 Since we have no channel-mixing, without loss of generality we can assume that  $d = 1$  – one can  
815 then extend the argument to any entry. We can use the Jordan form of  $\mathcal{A}$  to write the solution of the  
816 GRAND dynamical system as

$$\mathbf{f}(t) = P \text{diag}(e^{J_1 t}, \dots, e^{J_n t}) P^{-1} \mathbf{f}(0),$$

817 for some invertible matrix  $P$  of eigenvectors, with

$$e^{J_k t} = e^{-(1-\lambda_k^{\mathcal{A}})t} \begin{pmatrix} 1 & t & \dots & \frac{t^{m_k-1}}{(m_k-1)!} \\ & & \ddots & \\ & & & 1 \end{pmatrix},$$

818 where  $m_k$  are the eigenvalue multiplicities. Since by assumption  $G$  is connected and augmented with  
819 self-loops, the row-stochastic attention matrix  $\mathcal{A}$  computed in [10] with softmax activation is *regular*,  
820 meaning that there exists  $m \in \mathbb{N}$  such that  $(\mathcal{A}^m)_{ij} > 0$  for each entry  $(i, j)$ . Accordingly, we can  
821 apply Perron Theorem to derive that any eigenvalue of  $\mathcal{A}$  has real part smaller than one except the  
822 eigenvalue  $\lambda_0^{\mathcal{A}}$  with multiplicity one, associated with the Perron eigenvector  $\mathbf{1}_n$ . Accordingly, we  
823 find that each block  $e^{J_k t}$  decays to zero as  $t \rightarrow \infty$  with the exception of the one  $e^{J_0 t}$  associated with  
824 the Perron eigenvector. In particular, the projection of  $\mathbf{f}_0$  over the Perron eigenvector is just  $\mu \mathbf{1}_n$ , with  
825  $\mu$  the average of the feature initial condition. This completes the proof.  $\square$

## 826 C.2 Common GNN architectures as gradient flow

827 We consider linear GNNs of the form

$$\mathbf{F}(t+1) = \mathbf{F}(t) \Omega + \mathcal{A} \mathbf{F}(t) \mathbf{W} + \beta \mathbf{F}(0) \tilde{\mathbf{W}}, \quad 0 \leq t \leq T.$$

828 If  $\Omega = \mathbf{0}, \beta = 0$  and  $\mathcal{A} = \bar{\mathbf{A}}$ , we recover linear GCN with weights shared across layers [27, 43].  
829 Similarly, if  $\mathcal{A} = \bar{\mathbf{A}}$  and  $\beta = 0$ , this is linear GraphSAGE [23] with propagation given by *symmetric*  
830 adjacency and weights shared across layers. A symmetric version of GAT [42] can be recovered if  
831  $\Omega = \mathbf{0}, \beta = 0$  and  $\mathcal{A} = \bar{\mathbf{A}}$  is a *symmetric* attention matrix depending only on the initial encoded  
832 features – note that in general a row-stochastic matrix may not be symmetric so a *symmetrization*  
833 *of a row-stochastic attention matrix would generally fail to remain row-stochastic*. We believe that  
834 this point deserves further investigation. Finally GCNII [11] can be recovered by taking  $\Omega = \mathbf{0}$  and  
835  $\mathcal{A} = \bar{\mathbf{A}}$ .

836 *Proof of Lemma 4.2.* This follows from the same argument in eq. (22) once we regard the linear  
 837 system in eq. (12) as a unit step size Euler discretization

$$\dot{\mathbf{F}}(t) \sim \mathbf{F}(t+1) - \mathbf{F}(t) = \mathbf{F}(t)(\mathbf{\Omega} - \mathbf{I}_d) + \mathbf{A}\mathbf{F}(t)\mathbf{W} + \beta\mathbf{F}(0)\tilde{\mathbf{W}}$$

838 □

### 839 C.3 Spectral analysis of the channel-mixing: the discrete case

840 We first address the proof of the main result.

841 *Proof of Theorem 4.3.* We consider a linear dynamical system

$$\mathbf{F}(t + \tau) = \mathbf{F}(t) + \tau\bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W},$$

842 with  $\mathbf{W}$  symmetric. We vectorize the system and rewrite it as

$$\text{vec}(\mathbf{F}(t + \tau)) = (\mathbf{I}_{nd} + \tau\mathbf{W} \otimes \bar{\mathbf{A}})\text{vec}(\mathbf{F}(t))$$

843 which in particular leads to

$$\text{vec}(\mathbf{F}(m\tau)) = (\mathbf{I}_{nd} + \tau\mathbf{W} \otimes \bar{\mathbf{A}})^m \text{vec}(\mathbf{F}(0)).$$

844 We can then write explicitly the solution as

$$\text{vec}(\mathbf{F}(m\tau)) = \sum_{r,i} (1 + \tau\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}))^m c_{r,i}(0)\phi_r^{\mathbf{W}} \otimes \phi_i^{\Delta}.$$

845 We now verify that by assumption in eq. (13) the dominant term of the solution is the projection into  
 846 the eigenspace associated with the eigenvalue  $\rho_- = |\lambda_-^{\mathbf{W}}|(\rho_{\Delta} - 1)$ . The following argument follows  
 847 the same structure in the proof of Proposition 3.1 with the extra condition given by the step-size.  
 848 First, we note that for any  $r$  such that  $\lambda_r^{\mathbf{W}} > 0$ , we have

$$|1 + \tau\rho_-| > |1 + \tau\lambda_+^{\mathbf{W}}| \geq |1 + \tau\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta})|$$

849 since we required  $\rho_- > \lambda_+^{\mathbf{W}}$  in eq. (13). Conversely, if  $\lambda_r^{\mathbf{W}} < 0$ , then

$$|1 + \tau\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta})| \leq \max\{|1 + \tau\rho_-|, |1 + \tau\lambda_-^{\mathbf{W}}|\}$$

850 Assume that  $\tau|\lambda_-^{\mathbf{W}}| > 1$ , otherwise there is nothing to prove. Then  $|1 + \tau\rho_-| > \tau|\lambda_-^{\mathbf{W}}| - 1$  if and  
 851 only if

$$\tau|\lambda_-^{\mathbf{W}}|(2 - \rho_{\Delta}) < 2,$$

852 which is precisely the right inequality in eq. (13). We can then argue exactly as in the proof of  
 853 Proposition 3.1 to derive that for each index  $r$  such that  $\lambda_r^{\mathbf{W}} < 0$  and  $\lambda_r^{\mathbf{W}} \neq \lambda_-^{\mathbf{W}}$ , then

$$|1 + \tau\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta})| \leq \max\{|1 + \tau|\lambda_{-,2}^{\mathbf{W}}|(\rho_{\Delta} - 1)|, |1 + \tau|\lambda_-^{\mathbf{W}}|(\lambda_{n-2}^{\Delta} - 1)|\}$$

854 with  $\lambda_{-,2}^{\mathbf{W}}$  and  $\lambda_{n-2}^{\Delta}$  defined in eq. (23). We can then introduce

$$\delta_{\text{HFD}} := \max\{\lambda_+^{\mathbf{W}}, \rho_- - |\lambda_-^{\mathbf{W}}|\text{gap}(\rho_{\Delta}\mathbf{I} - \mathbf{\Delta}), \rho_- - (\rho_{\Delta} - 1)\text{gap}(|\lambda_-^{\mathbf{W}}|\mathbf{I} + \mathbf{W}), |\lambda_-^{\mathbf{W}}| - \frac{2}{\tau}\} \quad (28)$$

855 and conclude that

$$\begin{aligned} \mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) &= \frac{1}{2} \sum_{r,i} (1 + \tau\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}))^{2m} c_{r,i}^2(0)\lambda_i^{\Delta} \\ &= (1 + \tau\rho_-)^{2m} \left( \frac{\rho_{\Delta}}{2} \sum_{r:\lambda_r^{\mathbf{W}}=\lambda_-^{\mathbf{W}}} c_{r,\rho_{\Delta}}^2(0) + \mathcal{O}\left(\left(\frac{1 + \tau\delta_{\text{HFD}}}{1 + \tau\rho_-}\right)^{2m}\right) \sum_{i:r:\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}) \neq \rho_-} c_{r,i}^2(0)\lambda_i^{\Delta} \right) \\ &= (1 + \tau\rho_-)^{2m} \left( \frac{\rho_{\Delta}}{2} \|P_{\mathbf{W}}^{\rho_-} \mathbf{F}(0)\|^2 + \mathcal{O}\left(\left(\frac{1 + \tau\delta_{\text{HFD}}}{1 + \tau\rho_-}\right)^{2m}\right) \right). \end{aligned}$$



856 In particular, we can normalize the solution and due to  $(\mathbf{I}_d \otimes \Delta)P_{\mathbf{W}}^{\rho_-} \text{vec}(\mathbf{F}(0)) = \rho_{\Delta} P_{\mathbf{W}}^{\rho_-} \text{vec}(\mathbf{F}(0))$ ,  
 857 we complete the proof for the case with residual connection.

858 If instead we drop the residual connection and simply consider  $\dot{\mathbf{F}}(t) = \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$ , then

$$\text{vec}(\mathbf{F}(m\tau)) = (\tau\mathbf{W} \otimes \bar{\mathbf{A}})^m \text{vec}(\mathbf{F}(0)).$$

859 Since  $G$  is not bipartite, the Laplacian spectral radius satisfies  $\rho_{\Delta} < 2$ . Therefore, for each pair of  
 860 indices  $(r, i)$  we have the following bound:

$$|\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta})| \leq \max\{\lambda_+^{\mathbf{W}}, |\lambda_-^{\mathbf{W}}|\},$$

861 and the inequality becomes strict if  $i > 0$ , i.e.  $\lambda_i^{\Delta} > 0$ . The eigenvalues  $\lambda_+^{\mathbf{W}}$  and  $\lambda_-^{\mathbf{W}}$  are attained  
 862 along the eigenvectors  $\phi_+^{\mathbf{W}} \otimes \phi_0^{\Delta}$  and  $\phi_-^{\mathbf{W}} \otimes \phi_0^{\Delta}$  respectively. Accordingly, the dominant terms of the  
 863 evolution lie in the kernel of  $\mathbf{I}_d \otimes \Delta$ , meaning that for any  $\mathbf{F}_0$  with non-zero projection in  $\ker(\mathbf{I}_d \otimes \Delta)$   
 864 – which is satisfied by all initial conditions except those belonging to a lower dimensional subspace –  
 865 the dynamics is LFD. In fact, without loss of generality assume that  $|\lambda_-^{\mathbf{W}}| > \lambda_+^{\mathbf{W}}$ , then

$$\begin{aligned} \text{vec}(\mathbf{F}(m\tau)) &= |\lambda_-^{\mathbf{W}}|^m \sum_{r: \lambda_r^{\mathbf{W}} = \lambda_-^{\mathbf{W}}} (-1)^m c_{r,0}(0) \phi_-^{\mathbf{W}} \otimes \phi_0^{\Delta} \\ &\quad + |\lambda_-^{\mathbf{W}}|^m \left( \mathcal{O}(\varphi(m)) \left( \mathbf{I}_{nd} - \sum_{r: \lambda_r^{\mathbf{W}} = \lambda_-^{\mathbf{W}}} (\phi_-^{\mathbf{W}} \otimes \phi_0^{\Delta})(\phi_-^{\mathbf{W}} \otimes \phi_0^{\Delta})^{\top} \right) \text{vec}(\mathbf{F}(0)) \right), \end{aligned}$$

866 with  $\varphi(m) \rightarrow 0$  as  $m \rightarrow \infty$ , which completes the proof.  $\square$

867 **Gradient flow as spectral GNNs.** We finally discuss eq. (11) from the perspective of spectral  
 868 GNNs as in [2]. Let us assume that  $\beta = 0$ ,  $\Omega = \mathbf{0}$ . If we let  $\Delta = \mathbf{U}\Lambda\mathbf{U}^{\top}$  be the eigendecomposition  
 869 of the graph Laplacian and  $\{\lambda_r^{\mathbf{W}}\}$  be the spectrum of  $\mathbf{W}$  with associated orthonormal basis of  
 870 eigenvectors given by  $\{\phi_r^{\mathbf{W}}\}$ , and we introduce  $\mathbf{z}^r(t) : \mathcal{V} \rightarrow \mathbb{R}$  defined by  $z_i^r(t) = \langle \mathbf{f}_i(t), \phi_r^{\mathbf{W}} \rangle$ , then  
 871 we can rewrite the discretized gradient flow as

$$\mathbf{z}^r(t + \tau) = \mathbf{U}(\mathbf{I} + \tau\lambda_r^{\mathbf{W}}(\mathbf{I} - \Lambda))\mathbf{U}^{\top}\mathbf{z}^r(t) = \mathbf{z}^r(t) + \tau\lambda_r^{\mathbf{W}}\bar{\mathbf{A}}\mathbf{z}^r(t), \quad 1 \leq r \leq d. \quad (29)$$

872 Accordingly, for each projection into the  $r$ -th eigenvector of  $\mathbf{W}$ , we have a spectral function in the  
 873 graph frequency domain given by  $\lambda^{\Delta} \mapsto 1 + \tau\lambda_r^{\mathbf{W}}(1 - \lambda^{\Delta})$ . If  $\lambda_r^{\mathbf{W}} > 0$  we have a *low-pass* filter  
 874 while if  $\lambda_r^{\mathbf{W}} < 0$  we have a *high-pass* filter. Moreover, we see that along the eigenvectors of  $\mathbf{W}$ ,  
 875 if  $\lambda_r^{\mathbf{W}} < 0$  then the dynamics is equivalent to flipping the sign of the edge weights, which offers a  
 876 direct comparison with methods proposed in [4, 45] where some ‘attentive’ mechanism is proposed  
 877 to learn negative edge weights based on feature information.

878 The previous equation simply follows from

$$\begin{aligned} z_i^r(t + \tau) &= \langle \mathbf{f}_i(t + \tau), \phi_r^{\mathbf{W}} \rangle = \langle \mathbf{f}_i(t) + \mathbf{W}(\bar{\mathbf{A}}\mathbf{f}(t))_i, \phi_r^{\mathbf{W}} \rangle \\ &= z_i^r(t) + \lambda_r^{\mathbf{W}} \sum_j \bar{a}_{ij} z_j^r(t), \end{aligned}$$

879 which concludes the derivation of eq. (29).

## 880 D Additional details on experiments

### 881 D.1 Additional details on GRAFF

882 Given a gradient flow dynamical system of the form  $\mathbf{F}(t + \tau) = \mathbf{F}(t) + \tau\bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$ , the vectorized  
 883 solution is

$$\text{vec}(\mathbf{F}(m\tau)) = \sum_{r,i} (1 + \tau\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\Delta}))^m c_{r,i}(0) \phi_r^{\mathbf{W}} \otimes \phi_i^{\Delta}.$$

884 We then see that the number of layers  $m$  – which coincides with the quotient of the integration  
 885 time  $T$  by the step size  $\tau$  – represents the degree of the polynomial computing the solution. More

precisely, on a heterophilic graph for which a HFD dynamics is more suited than an LFD dynamics, the negative eigenvalues of  $\mathbf{W}$  are needed to magnify the graph high-frequencies. This in turn yields terms  $(1 + \tau\lambda_r^{\mathbf{W}}(1 - \lambda_i^{\mathbf{A}}))^m$  that would become unbounded with  $m$  growing if there is sufficient mass on the negative side of  $\text{spec}(\mathbf{W})$ . On the other hand, terms associated with positive eigenvalues of  $\mathbf{W}$  would quickly lead to over-smoothing. One then expects that on heterophilic graphs the degree  $m$  of the polynomial – i.e. the *number of layers* – should be generally smaller than that on homophilic graphs. This is confirmed in our real-world experiments where on the larger heterophilic graphs like Squirrel and Chameleon the optimal number  $m$  is an integer in  $\{2, 3, 4\}$ .

## D.2 General Experimental details

GRAFF is implemented in PyTorch [53], using PyTorch geometric [54] and torchdiffeq [12]. Code and instructions to reproduce the experiments are available on GitHub. Hyperparameters were tuned using wandb[55] and random grid search. Experiments were run on AWS p2.xlarge machines, each with 8 Tesla V100-SXM2 GPUs.

## D.3 Additional details on synthetic ablation studies:

The synthetic Cora dataset is provided by [51, Appendix G]. They use a modified preferential attachment process to generate graphs for target levels of homophily. Nodes, edges and features are sampled from Cora proportional to a mix of class compatibility and node degree resulting in a graph with the required homophily and appropriate feature/label distribution. To validate the provided data before use we provide table 2 summarising the properties of the synthetic Cora dataset. All rows/levels of homophily have the same number of nodes (1,490), edges (5,936), features (1,433) and classes (5).

homophily	max_degree	min_degree	av_degree	density	edge_homoph	node_homoph
0.00	84.33	1.67	3.98	0.0027	0.00	0.00
0.10	71.33	2.00	3.98	0.0027	0.10	0.10
0.20	73.33	1.67	3.98	0.0027	0.20	0.20
0.30	70.00	2.00	3.98	0.0027	0.29	0.30
0.40	77.67	2.00	3.98	0.0027	0.39	0.39
0.50	76.33	2.00	3.98	0.0027	0.49	0.49
0.60	76.00	1.67	3.98	0.0027	0.59	0.60
0.70	67.67	2.00	3.98	0.0027	0.70	0.70
0.80	58.00	1.67	3.98	0.0027	0.78	0.79
0.90	58.00	1.67	3.98	0.0027	0.89	0.89
1.00	51.00	2.00	3.98	0.0027	1.00	1.00

Table 2: Summary of properties of synthetic Cora dataset

As well as the ablation shown in fig. 2 we used this dataset to perform an ablation using GCN as the baseline. We asses the impact of each of the steps necessary to augment a standard GCN model to GRAFF. This involves 5 steps; 1) add an encoder/decoder. 2) add a residual connection. 3) share the weights of  $\mathbf{W}$  and  $\mathbf{\Omega}$  across time/layers. 4) symmetrize  $\mathbf{W}$  and  $\mathbf{\Omega}$ . 5) remove the non-linearity between layers. The results are shown in fig. 3 and corroborate Theorem 4.3 that adding a residual term is beneficial especially in low-homophily scenarios. We also note augmentations 3,4 and 5 are not "costly" in terms of performance.

## D.4 Additional details on real-world ablation studies

For the real-world experiments in table 1 we performed 10 repetitions over the splits taken from [31]. For all datasets we used the largest connected component (LCC) apart from Citeseer where the 5th and 6th split are LCC and others require the full dataset. For Chameleon and Squirrel we added self loops and made the edges undirected as a preprocessing step. All other datasets are provided as undirected but without self loops. Each split uses 48/32/20 of nodes for training, validation and test set respectively. Table 3 summarises each of the datasets.

We used the real-world datasets to perform 2 ablation studies. First we choose 2 heterophilic datasets (Chameleon, Squirrel) and 2 homophilic (Cora, Citeseer) and observed how the size of the hidden

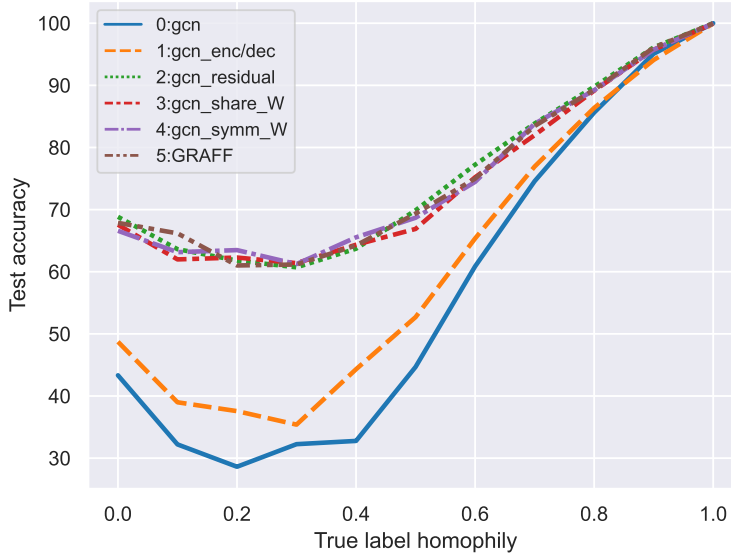


Figure 3: Experiments on synthetic Cora - GCN ablation

dataset	nodes	edges	features	classes	max_degree	min_degree	av_degree	density	edge_homoph	node_homoph
Texas	183	558	1,703	5	104	1	3.05	0.0167	0.06	0.06
Wisconsin	251	900	1,703	5	122	1	3.59	0.0143	0.18	0.16
Cornell	183	554	1,703	5	94	1	3.03	0.0165	0.3	0.3
Film	7,600	53,318	932	5	1,303	1	7.02	0.0009	0.22	0.22
Squirrel	5,201	401,907	2,089	5	1,904	2	77.27	0.0149	0.23	0.29
Chameleon	2,277	65,019	2,325	5	733	2	28.55	0.0125	0.26	0.33
Citeseer *	3,327	9,104	3,703	6	99	0	2.74	0.0008	0.74	0.71
Citeseer	2,120	7,358	3,703	6	99	1	3.47	0.0016	0.73	0.71
Pubmed	19,717	88,648	500	3	171	1	4.5	0.0002	0.8	0.79
Cora	2,485	10,138	1,433	7	168	1	4.08	0.0016	0.8	0.81

Table 3: Summary of properties of real-world datasets. All LCC except \*

dimension effected performance for the structures of  $\mathbf{W}$  described in section 5. For heterophilic datasets we used the splits from [31]. For homophilic datasets we used the methodology in [52], each split randomly selects 1,500 nodes for the development set, from the development set 20 nodes for each class are taken as the training set, the remainder are allocated as the validation set. The remaining nodes outside of the development set are used as the test set. This gives a lower percentage (3-6%) of training nodes. This approach was taken because less training information is needed in the homophilic setting and performance can become less sensitive to other factors, meaning less signal from the controlled variable. From fig. 4 we see that (DD) is more parameter efficient than *sum* in the heterophilic setting and (D) (a parameter light configuration) outperforms in the homophilic setting.

The second ablation study further corroborates the behaviour seen in fig. 2. We tested the structures of  $\mathbf{W}$  against the real-world datasets with known homophily, again *neg-prod* outperforms *prod* in the heterophilic setting and vice-versa due the sign of their spectra.

dataset	neg_prod	prod	sum
Chameleon	67.32	58.86	68.36
Squirrel	51.39	42.11	51.29
Cora	31.80	79.65	81.17
Citeseer	32.47	67.31	67.53

Table 4: Ablation with controlled spectrum of  $\mathbf{W}$  on real-world datasets

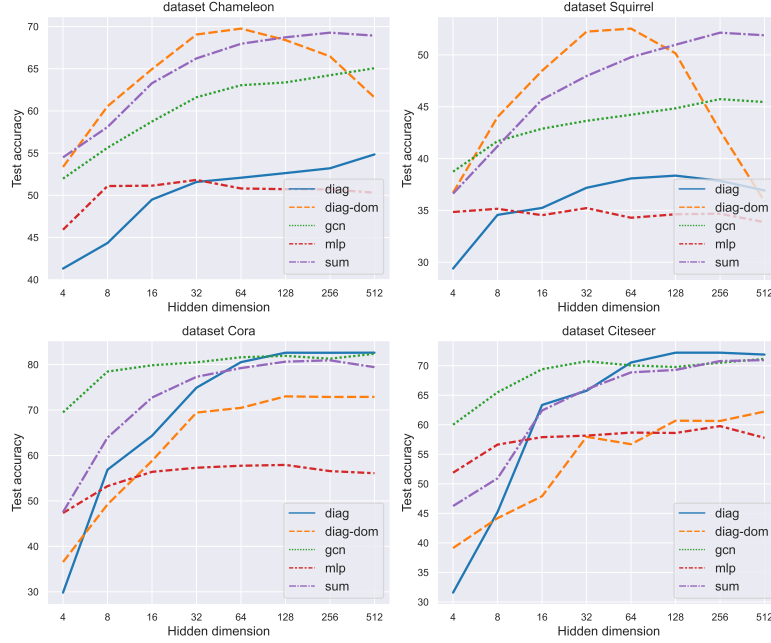


Figure 4: Ablation against hidden dimension

To validate the complexity analysis in Section 5 we performed a runtime ablation for the models between standard GCN and GRAFF described in the GCN ablation Figure 3. The average inference runtime over 100 runs for 1 split of Cora was recorded. We also include runtimes for the provided dense and sparse implementations of GGCN [45]. Adding the encoder/decoder (step 1) speeds up the model due to dimensionality reduction. Subsequent steps also reduce complexity and offer speedup with GRAFF performing the fastest.

## D.5 Details on hyperparameters

Using wandb [55] we performed a random grid search with uniform sampling of the continuous variables. We provide the hyperparameters that achieved the best results from the random grid search in table 5. An implementation that uses these hyperparameters is available in the provided code with hyperparameters provided in `graff_params.py`. Input dropout and dropout are the rates applied to the encoder/decoder respectively *with no dropout applied in the ODE block*.

	w_style	lr	decay	dropout	input_dropout	hidden_dim	time	step_size
chameleon	diag_dom	0.0014	0.0004	0.37	0.43	64	3.2	1
squirrel	diag_dom	0.0058	0.0002	0.50	0.51	64	2.3	1
texas	diag_dom	0.0041	0.0354	0.33	0.39	64	0.6	0.5
wisconsin	diag	0.0029	0.0318	0.37	0.37	64	2.1	0.5
cornell	diag	0.0021	0.0184	0.30	0.44	64	2.0	1
film	diag	0.0026	0.0130	0.48	0.42	64	1.5	1
Cora	diag	0.0026	0.0413	0.34	0.53	64	3.0	0.25
Citeseer	diag	0.0001	0.0274	0.22	0.51	64	2.0	0.5
Pubmed	diag	0.0039	0.0003	0.42	0.41	64	2.6	0.5

Table 5: Selected hyperparameters for real-world datasets

## E Elementwise non-linear activations and energy dissipation

In this section we investigate how to extend the energy framework to include more conventional non-linear activation maps to potentially have better expressive power – note that recent works like

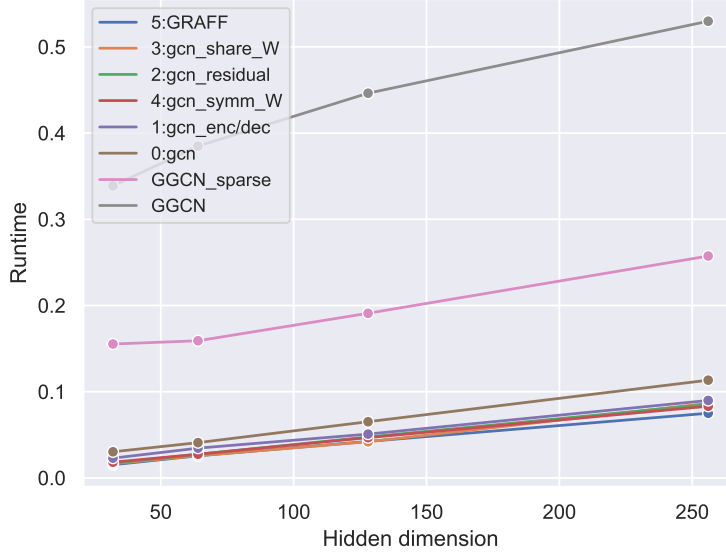


Figure 5: Runtime ablation for inference on Cora dataset

[6] in some sense argue for more non-linear layers. Namely, we consider the general energy in equation 7 with the inclusion of the source term:

$$\mathcal{E}^{\text{tot}}(\mathbf{F}) = \frac{1}{2} \sum_i \langle \mathbf{f}_i, \mathbf{\Omega} \mathbf{f}_i \rangle - \frac{1}{2} \sum_{i,j} \bar{a}_{ij} \langle \mathbf{f}_i, \mathbf{W} \mathbf{f}_j \rangle + \beta \langle \mathbf{F}, \mathbf{F}(0) \rangle.$$

Since the energy is *quadratic*, its gradient flow is a linear dynamical system:

$$\dot{\mathbf{F}}(t) = -\mathbf{F}(t)\mathbf{\Omega} + \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W} - \beta\mathbf{F}(0) \quad (30)$$

The key question we explore here is: **what happens if we activate the equations with a pointwise non-linear map  $\sigma$** ? In general, we will not be a gradient flow of the energy  $\mathcal{E}^{\text{tot}}$ , however can we still say something about the behaviour of  $t \mapsto \mathcal{E}^{\text{tot}}(\mathbf{F}(t))$  along the solution? The answer is affirmative and offers a novel contribution where even non-linear, residual, graph convolutional models maintain the interpretation of dissipating an energy where  $\mathbf{W}$  plays the role of an edge-wise bilinear potential generating attraction and repulsion:

**Proposition E.1.** Consider a non-linear map  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  such that the function  $x \mapsto x\sigma(x) \geq 0$ . If  $t \mapsto \mathbf{F}(t)$  solves the equation

$$\dot{\mathbf{F}}(t) = \sigma \left( -\mathbf{F}(t)\mathbf{\Omega} + \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W} - \beta\mathbf{F}(0) \right),$$

where  $\sigma$  acts elementwise, then

$$\frac{d\mathcal{E}^{\text{tot}}(\mathbf{F}(t))}{dt} \leq 0.$$

*Proof.* The argument is simple and derives from direct computation. Namely, let us use the Kronecker product formalism to rewrite the gradient  $\nabla_{\text{vec}(\mathbf{F})}\mathcal{E}^{\text{tot}}(\mathbf{F})$  as a vector in  $\mathbb{R}^{nd}$ : explicitly, we get

$$\nabla_{\text{vec}(\mathbf{F})}\mathcal{E}^{\text{tot}}(\mathbf{F}) = (\mathbf{\Omega} \otimes \mathbf{I}_n - \mathbf{W} \otimes \bar{\mathbf{A}})\text{vec}(\mathbf{F}) + \beta\text{vec}(\mathbf{F}(0)).$$

It follows then that

$$\begin{aligned} \frac{d\mathcal{E}^{\text{tot}}(\mathbf{F}(t))}{dt} &= (\nabla_{\text{vec}(\mathbf{F})}\mathcal{E}^{\text{tot}}(\mathbf{F}(t)))^\top \text{vec}(\dot{\mathbf{F}}(t)) = \\ &= ((\mathbf{\Omega} \otimes \mathbf{I}_n - \mathbf{W} \otimes \bar{\mathbf{A}})\text{vec}(\mathbf{F}(t)) + \beta\text{vec}(\mathbf{F}(0)))^\top \sigma \left( (-\mathbf{\Omega} \otimes \mathbf{I}_n + \mathbf{W} \otimes \bar{\mathbf{A}})\text{vec}(\mathbf{F}(t)) - \beta\text{vec}(\mathbf{F}(0)) \right) \end{aligned}$$

965 If we introduce the notation  $\mathbf{Z}(t) = (-\mathbf{\Omega} \otimes \mathbf{I}_n + \mathbf{W} \otimes \bar{\mathbf{A}})\text{vec}(\mathbf{F}(t)) - \beta\text{vec}(\mathbf{F}(0))$ , then we can  
 966 rewrite the derivative as

$$\frac{d\mathcal{E}^{\text{tot}}(\mathbf{F}(t))}{dt} = -\mathbf{Z}(t)^\top \sigma(\mathbf{Z}(t)) = -\sum_{\alpha} \mathbf{Z}(t)^\alpha \sigma(\mathbf{Z}(t)^\alpha) \leq 0$$

967 by assumption on  $\sigma$ , which completes the proof.  $\square$

968 **Important consequence:** The previous results shows that even if the non-linear dynamical system

$$\dot{\mathbf{F}}(t) = \sigma(-\mathbf{F}(t)\mathbf{\Omega} + \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W} + \beta\mathbf{F}(0)),$$

969 is not a gradient flow for  $\mathcal{E}^{\text{tot}}$ , the latter quantity is still decreasing along the solution meaning that  
 970 the interpretation of positive (negative) eigenvalues of  $\mathbf{W}$  inducing attraction (repulsion) persists  
 971 given that *the energy has not changed*. **This allows us to derive that general (non-linear) graph**  
 972 **convolutional models retain the learnable multi-particle energy-dissipation property provided**  
 973 **that the channel-mixing matrices are symmetric and that the pointwise activation satisfies**  
 974  $x\sigma(x) \geq 0$ , which for example holds for ReLU, tanh, arctan and so on. *In particular, models that*  
 975 *are energy-dissipating can fit non-linear activations.*

976 To further support the principle that the effects induced by  $\mathbf{W}$  are similar even in this non-linear  
 977 setting, we consider a simplified scenario.

978 **Lemma E.2.** *If we choose  $\mathbf{\Omega} = \mathbf{W} = \text{diag}(\omega)$  with  $\omega^r \leq 0$  for  $1 \leq r \leq d$  and  $\beta = 0$  i.e.  $t \mapsto \mathbf{F}(t)$*   
 979 *solves the dynamical system*

$$\dot{\mathbf{F}}(t) = \sigma(-\mathbf{\Delta}\mathbf{F}(t)\text{diag}(\omega)),$$

980 *with  $x\sigma(x) \geq 0$ , then the standard graph Dirichlet energy satisfies*

$$\frac{d\mathcal{E}^{\text{Dir}}(\mathbf{F}(t))}{dt} \geq 0.$$

981 *Proof.* This again simply follows from directly computing the derivative:

$$\begin{aligned} \frac{d\mathcal{E}^{\text{Dir}}(\mathbf{F}(t))}{dt} &= \frac{1}{4} \frac{d}{dt} \left( \sum_{r=1}^d \sum_{(i,j) \in \mathbf{E}} \left( \frac{f_i^r(t)}{\sqrt{d_i}} - \frac{f_j^r(t)}{\sqrt{d_j}} \right)^2 \right) \\ &= \sum_{r=1}^d \sum_{i \in \mathbf{V}} (\mathbf{\Delta}\mathbf{f}^r)_i \sigma(-\omega^r (\mathbf{\Delta}\mathbf{f}^r)_i) = \sum_{r=1}^d \sum_{i \in \mathbf{V}} (\mathbf{\Delta}\mathbf{f}^r)_i \sigma(|\omega^r| (\mathbf{\Delta}\mathbf{f}^r)_i) \geq 0. \end{aligned}$$

982 **Important consequence:** The previous Lemma implies that even with non-linear activations, negative  
 983 eigenvalues of the channel-mixing induce repulsion and indeed the solution becomes less smooth as  
 984 measured by the classical Dirichlet Energy increasing along the solution. Generalising this result to  
 985 more arbitrary choices is not immediate and we reserve this for future work.

986  $\square$