
Graph Neural Networks as Gradient Flows

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

13 1 Introduction and motivations

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

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

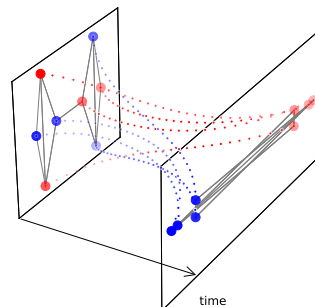


Figure 1: GRAFF dynamics: attractive and repulsive forces lead to a non-smoothing process able to separate labels.

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

45 **Related work.** Our analysis is related to studying GNNs as filters on the graph spectrum [15, 24,
46 2, 25] and over-smoothing [29, 30, 8, 50] and partly adopts techniques similar to [30]. The key
47 difference is that we also consider the spectrum of the ‘channel-mixing’ matrix. The concept of
48 gradient flows has been a standard tool in physics and geometry [16], from which they were adopted
49 for image processing [26], and recently used in ML [35] for the analysis of Transformers [41] – see
50 also [18] for discussion of loss landscapes. Our continuous-time evolution equations follows the spirit
51 of Neural ODES [22, 12, 3] and the study of GNNs as continuous dynamical systems [44, 10, 17, 9].

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

67 2 Gradient-flow formalism

68 **Notations adopted throughout the paper.** Let $G = (V, E)$ be an *undirected* graph with n nodes.
69 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
70 (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
71 by stacking its columns. Given a symmetric matrix \mathbf{B} , we let $\lambda_+^{\mathbf{B}}, \lambda_-^{\mathbf{B}}$ denote its most positive and
72 negative eigenvalues, respectively, and $\rho_{\mathbf{B}}$ be its *spectral radius*. If $\mathbf{B} \succeq 0$, then $\text{gap}(\mathbf{B})$ denotes the
73 *positive smallest eigenvalue* of \mathbf{B} . $\dot{f}(t)$ denotes the temporal derivative, \otimes is the Kronecker product
74 and ‘a.e.’ means *almost every* w.r.t. Lebesgue measure and usually refers to data in the complement
75 of some lower dimensional subspace in $\mathbb{R}^{n \times d}$. Proofs and additional results appear in the Appendix.

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

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

79 and measures the ‘smoothness’ of f . A natural approach to find minimizers of \mathcal{E} - called *harmonic*
80 *maps* - was introduced in [16] and consists in studying the **gradient flow** of \mathcal{E} , wherein a given map
81 $f(0) = f_0$ is evolved according to $\dot{f}(t) = -\nabla_f \mathcal{E}(f(t))$. These type of evolution equations have
82 historically been the core of *variational* and *PDE-based image processing*; in particular, gradient
83 flows of the Dirichlet energy were shown [26] to recover the Perona-Malik nonlinear diffusion [32].

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

87 parametric rule, $\text{GNN}_\theta(\mathbf{G}, \mathbf{F}_0)$, and uses a decoding map for the prediction $y = \psi_{\text{DE}}(\text{GNN}_\theta(\mathbf{G}, \mathbf{F}_0))$.
 88 In graph convolutional models [15][27], GNN_θ consists of two operations: applying a shared linear
 89 transformation to the features (**‘channel mixing’**) and propagating them along the edges of the graph
 90 (**‘diffusion’**). Our **goal** consists in studying when GNN_θ is the *gradient flow* of some parametric class
 91 of energies $\mathcal{E}_\theta : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}$, which generalize the Dirichlet energy. This means that the parameters
 92 can be interpreted as ‘finding the right notion of smoothness’ for our task. We evolve the features by
 93 $\dot{\mathbf{F}}(t) = -\nabla_{\mathbf{F}} \mathcal{E}_\theta(\mathbf{F}(t))$ with prediction $y = \psi_{\text{DE}}(\mathbf{F}(T))$ for some optimal time T .

94 **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
 95 flow. Accordingly, this framework allows to *explain the GNN dynamics* as flowing the node features
 96 in the direction of steepest descent of \mathcal{E}_θ . Indeed, we find that parametrizing an energy leads to
 97 equations governed by attractive and repulsive forces that can be controlled via the spectrum of
 98 symmetric ‘channel-mixing’ matrices. This shows that by learning to distribute more mass over the
 99 negative (positive) eigenvalues of the channel-mixing, gradient flow models can generate dynamics
 100 dominated by the higher (respectively, lower) graph frequencies and hence tackle different homophily
 101 scenarios. The gradient flow framework also leads to sharing of the weights across layers (since we
 102 parametrize the *energy* rather than the *evolution equations*, as usually done in GNNs), allowing us to
 103 reduce the number of parameters without compromising performance (see Table 1).

104 **Analysis on graphs: preliminaries.** Given a *connected* graph \mathbf{G} with self-loops, its adjacency
 105 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
 106 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
 107 *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
 108 *divergence* div is the adjoint of ∇), represented by $\Delta = \mathbf{I} - \bar{\mathbf{A}} \geq 0$. We refer to the eigenvalues of
 109 Δ as *frequencies*: the lowest frequency is always 0 while the highest frequency is $\rho_\Delta \leq 2$ [14]. As
 110 for the continuum case, the gradient allows to define a (*graph*) *Dirichlet energy* as [49]

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

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

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

113 for $1 \leq r \leq d$. Similarly to [8], we rely on \mathcal{E}^{Dir} to assess whether a given dynamics $t \mapsto \mathbf{F}(t)$ is a
 114 smoothing process. A different choice of Laplacian $\mathbf{L} = \mathbf{D} - \mathbf{A}$ with non-normalized adjacency
 115 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
 116 on the following definitions (see Appendix A.3 for further equivalent formulations and justifications):

117 **Definition 2.1.** $\dot{\mathbf{F}}(t) = \text{GNN}_\theta(\mathbf{F}(t), t)$ initialized at $\mathbf{F}(0)$ is *smoothing* if $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \leq C + \varphi(t)$,
 118 with C a constant only depending on $\mathcal{E}^{\text{Dir}}(\mathbf{F}(0))$ and $\varphi(t) \leq 0$. *Over-smoothing* occurs if either
 119 $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow 0$ or $\mathcal{E}_\mathbf{L}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow 0$ for $t \rightarrow \infty$.

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

125 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.
 126 $\mathbf{F}(0)$, the low-frequency components are growing the fastest and indeed $\mathbf{F}(t)/\|\mathbf{F}(t)\| \rightarrow \mathbf{F}_\infty$ s.t.
 127 $\Delta \mathbf{f}_\infty^r = \mathbf{0}$ for $1 \leq r \leq d$. We formalize this scenario – including the opposite case of high-frequency
 128 components being dominant – by studying $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)/\|\mathbf{F}(t)\|)$, i.e. the Rayleigh quotient of $\mathbf{I}_d \otimes \Delta$.

129 **Definition 2.2.** $\dot{\mathbf{F}}(t) = \text{GNN}_\theta(\mathbf{F}(t), t)$ initialized at $\mathbf{F}(0)$ is *Low/High-Frequency-Dominant*
 130 (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 $t \rightarrow \infty$.

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

133 **Lemma 2.3.** GNN_θ is LFD (HFD) iff for each $t_j \rightarrow \infty$ there exist $t_{j_k} \rightarrow \infty$ and \mathbf{F}_∞ s.t.
 134 $\mathbf{F}(t_{j_k})/\|\mathbf{F}(t_{j_k})\| \rightarrow \mathbf{F}_\infty$ 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 Δ associated with largest frequency ρ_Δ 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}^{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}^{Dir} as

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

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 $(i, j) \in E$.
 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 $i \in \mathcal{V}$ 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 Ω 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}^{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

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

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

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

176 **Spectral analysis of the channel-mixing.** We will now show that eq. (8) can lead to a HFD
 177 dynamics. To this end, we assume that $\Omega = \mathbf{0}$ so that eq. (8) becomes $\dot{\mathbf{F}}(t) = \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$. According
 178 to eq. (9) the negative eigenvalues of \mathbf{W} lead to repulsion. We show that the latter can induce HFD
 179 dynamics as per Definition 2.2. We let $P_{\mathbf{W}}^{\rho_-}$ be the orthogonal projection into the eigenspace of
 180 $\mathbf{W} \otimes \bar{\mathbf{A}}$ associated with the eigenvalue $\rho_- := |\lambda_-^{\mathbf{W}}|(\rho_\Delta - 1)$. We define ϵ_{HFD} explicitly in eq. (24).

181 **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 ϵ_{HFD} s.t.*

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

182 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 $1 \leq r \leq d$.

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

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

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

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

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

199 **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*
 200 *$t \mapsto \mathbf{F}(t)$ solves the equation*

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

201 *where σ acts elementwise, then*

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

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

206 4 Comparison with GNNs

207 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_D 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_D, CGNN and GRAND satisfy the following:

- 218 (i) PDE – GCN_D 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.

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

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

255 **Over-smoothing analysis in discrete setting.** By Proposition 3.1 we know that the continuous
 256 version of eq. (11) can be HFD thanks to the negative eigenvalues of \mathbf{W} . The next result represents a
 257 discrete counterpart of Proposition 3.1 and shows that *residual, symmetrized graph convolutional*
 258 *models can be HFD*. Below $P_{\mathbf{W}}^{\rho_-}$ is the projection into the eigenspace associated with the eigenvalue
 259 $\rho_- := |\lambda_-^{\mathbf{W}}|(\rho_{\Delta} - 1)$ and we report the explicit value of δ_{HFD} in eq. (28) in Appendix C.3. We let:

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

260 **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

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

261 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. $\Delta \mathbf{f}_{\infty}^r = \rho_{\Delta} \mathbf{f}_{\infty}^r$.
 262 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}
 263 symmetric, is LFD independent of the spectrum of \mathbf{W} .

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

271 **The residual connection from a spectral perspective.** Given a sufficiently small step-size so
 272 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.
 273 $\mathbf{F}(0)$ 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
 274 positive one. This means that differently from [29, 30, 8], there is no requirement on the minimal
 275 magnitude of the spectral radius of \mathbf{W} coming from the graph topology as long as $\lambda_+^{\mathbf{W}}$ is small
 276 enough. Conversely, without a residual term, the dynamics is LFD for a.e. $\mathbf{F}(0)$ independently of the
 277 sign and magnitude of the eigenvalues of \mathbf{W} . This is also confirmed by the GCN-curve in Figure 2.

278 **Over-smoothing vs LFD.** We highlight how in general a linear GCN equation as $\mathbf{F}(t + \tau) =$
 279 $\tau \bar{\mathbf{A}}\mathbf{F}(t)\mathbf{W}$ may avoid over-smoothing in the sense of Definition 2.1 meaning that $\mathcal{E}^{\text{Dir}}(\mathbf{F}(t)) \rightarrow \infty$
 280 as soon as there exist $\lambda_i^{\Delta} \in (0, 1)$ and the spectral radius of \mathbf{W} is large enough. However, this
 281 will not lead to over-separation since the dominating term is the lowest frequency one: in other
 282 words, once we re-set the scale right as per the normalization in Theorem 4.3 we encounter loss of
 283 separability even with large (and possibly negative) spectrum of \mathbf{W} .

284 5 Experiments

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

291 **Methodology.** We crystallize GRAFF in the model presented in eq. (11) with $\psi_{\text{EN}}, \psi_{\text{DE}}$ im-
 292 plemented as single linear layers or MLPs, and we set Ω to be diagonal. For the real-world
 293 experiments we consider *diagonally-dominant* (DD), *diagonal* (D) and *time-dependent* choices
 294 for the structure of \mathbf{W} that offer explicit control over its spectrum. In the (DD)-case, we consider
 295 a $\mathbf{W}^0 \in \mathbb{R}^{d \times d}$ symmetric with zero diagonal and $\mathbf{w} \in \mathbb{R}^d$ defined by $w_\alpha = q_\alpha \sum_\beta |\mathbf{W}_{\alpha\beta}^0| + r_\alpha$,
 296 and set $\mathbf{W} = \text{diag}(\mathbf{w}) + \mathbf{W}^0$. Due to the Gershgorin Theorem the eigenvalues of \mathbf{W} belong to
 297 $[\mathbf{w}_\alpha - \sum_\beta |\mathbf{W}_{\alpha\beta}^0|, \mathbf{w}_\alpha + \sum_\beta |\mathbf{W}_{\alpha\beta}^0|]$, so the model ‘can’ easily re-distribute mass in the spectrum of
 298 \mathbf{W} via q_α, r_α . This generalizes the decomposition of \mathbf{W} in (11) providing a justification in terms of
 299 its spectrum and turns out to be more efficient w.r.t. the hidden dimension d as shown in Figure 4
 300 in the Appendix. For (D) we take \mathbf{W} to be diagonal, with entries sampled $\mathcal{U}[-1, 1]$ and fixed – i.e., we
 301 **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
 302 varies across layers. To investigate the role of the spectrum of \mathbf{W} on synthetic graphs, we construct
 303 three additional variants: $\mathbf{W} = \mathbf{W}' + \mathbf{W}'^\top$, $\mathbf{W} = \pm \mathbf{W}'^\top \mathbf{W}'$ named *sum*, *prod* and *neg-prod*
 304 respectively where *prod* (*neg-prod*) variants have only non-negative (non-positive) eigenvalues.

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

313 Synthetic experiments and ablation studies.

314 To investigate our claims in a controlled environment we use the synthetic Cora dataset of [51] Appendix G].
 315 Graphs are generated for target levels of homophily via preferential attachment – see
 316 Appendix D.3 for details. Figure 2 confirms the spectral analysis and offers a better understanding
 317 in terms of performance and smoothness of the predictions. Each curve – except GCN – represents
 318 one version of \mathbf{W} as in ‘methodology’ and we implement eq. (11) with $\beta = 0, \Omega = \mathbf{0}$. Figure
 319 2 (top) reports the test accuracy vs true label homophily. *Neg-prod* is better than *prod* on low-
 320 homophily and viceversa on high-homophily. This confirms Proposition 3.1 where we have shown
 321 that the gradient flow can lead to a HFD dynamics – that are generally desirable with low-
 322 homophily – through the negative eigenvalues of \mathbf{W} . Conversely, the *prod* configuration (where we
 323 have an attraction-only dynamics) struggles in low-homophily scenarios *even though a residual connection is present*.
 324 Both *prod* and *neg-prod* are ‘extreme’ choices and serve the purpose of highlighting that by turning off one side of the spectrum
 325 this could be the more damaging depending on the underlying homophily. In general though ‘neutral’
 326 variants like *sum* and (DD) are indeed more flexible and better performing. In fact, (DD) outperforms
 327 GCN especially in low-homophily scenarios, confirming Theorem 4.3 where we have shown that
 328 without a residual connection convolutional models are LFD – and hence more sensitive to underlying
 329 homophily – irrespectively of the spectrum of \mathbf{W} . This is further confirmed in Figure 3

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

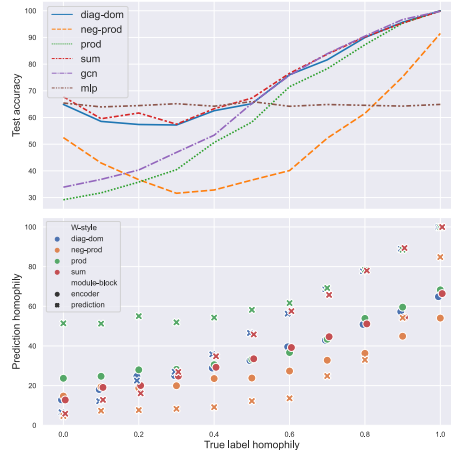


Figure 2: Experiments on synthetic datasets with controlled homophily.

	Texas	Wisconsin	Cornell	Film	Squirrel	Chameleon	Citeseer	Pubmed	Cora
Hom level	0.11	0.21	0.30	0.22	0.22	0.23	0.74	0.80	0.81
#Nodes	183	251	183	7,600	5,201	2,277	3,327	18,717	2,708
#Edges	295	466	280	26,752	198,493	31,421	4,676	44,327	5,278
#Classes	5	5	5	5	5	5	7	3	6
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

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

350 **Real world experiments.** We test GRAFF against a range of datasets with varying homophily
351 [37, 33, 31] (see Appendix D.4 for additional details). We use results provided in [45] Table 1],
352 which includes standard baselines as GCN [27], GraphSAGE [23], GAT [42], PairNorm [48] and
353 recent models tailored towards the heterophilic setting (GGCN [45], Geom-GCN [31], H2GCN
354 [51] and GPRGNN [13]). For Sheaf [5], a recent top-performer on heterophilic datasets, we took
355 the best performing variant (out of six provided) for each dataset. We also include continuous
356 baselines CGNN [44] and GRAND [10] to provide empirical evidence for Proposition 4.1. Splits
357 taken from [31] are used in all the comparisons. The GRAFF model discussed in ‘methodology’
358 is a very simple architecture with shared parameters across layers and run-time smaller than GCN
359 and more recent models like GGCN designed for heterophilic graphs (see Figure 5 in the Appendix).
360 Nevertheless, it achieves competitive results on all datasets, performing on par or better than more
361 complex recent models. Moreover, comparison with the ‘time-dependent’ (DD) variant confirms
362 that by sharing weights across layers we do not lose performance. We note that on heterophilic
363 graphs short integration time is usually needed due to the topology being harmful and the negative
364 eigenvalues of \mathbf{W} leading to exponential behaviour (see Appendix D).

365 6 Conclusions

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

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

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

386 **References**

- 387 [1] U. Alon and E. Yahav. On the bottleneck of graph neural networks and its practical implications.
388 In *International Conference on Learning Representations*, 2021.
- 389 [2] M. Balcilar, G. Renton, P. Héroux, B. Gaüzère, S. Adam, and P. Honeine. Analyzing the
390 expressive power of graph neural networks in a spectral perspective. In *International Conference*
391 *on Learning Representations*, 2020.
- 392 [3] M. Biloš, J. Sommer, S. S. Rangapuram, T. Januschowski, and S. Günnemann. Neural flows:
393 Efficient alternative to neural odes. In *Advances in Neural Information Processing Systems*,
394 volume 34, 2021.
- 395 [4] D. Bo, X. Wang, C. Shi, and H. Shen. Beyond low-frequency information in graph convolutional
396 networks. In *AAAI. AAAI Press*, 2021.
- 397 [5] C. Bodnar, F. Di Giovanni, B. P. Chamberlain, P. Liò, and M. M. Bronstein. Neural sheaf
398 diffusion: A topological perspective on heterophily and oversmoothing in gnns. *arXiv preprint*
399 *arXiv:2202.04579*, 2022.
- 400 [6] S. Brody, U. Alon, and E. Yahav. How attentive are graph attention networks? *arXiv preprint*
401 *arXiv:2105.14491*, 2021.
- 402 [7] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected
403 networks on graphs. In *2nd International Conference on Learning Representations, ICLR 2014*,
404 2014.
- 405 [8] C. Cai and Y. Wang. A note on over-smoothing for graph neural networks. *arXiv preprint*
406 *arXiv:2006.13318*, 2020.
- 407 [9] B. Chamberlain, J. Rowbottom, D. Eynard, F. Di Giovanni, X. Dong, and M. Bronstein. Beltrami
408 flow and neural diffusion on graphs. *Advances in Neural Information Processing Systems*, 34,
409 2021.
- 410 [10] B. Chamberlain, J. Rowbottom, M. I. Gorinova, M. Bronstein, S. Webb, and E. Rossi. Grand:
411 Graph neural diffusion. In *International Conference on Machine Learning*, pages 1407–1418.
412 PMLR, 2021.
- 413 [11] M. Chen, Z. Wei, Z. Huang, B. Ding, and Y. Li. Simple and deep graph convolutional networks.
414 In *International Conference on Machine Learning*, pages 1725–1735. PMLR, 2020.
- 415 [12] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural ordinary differential
416 equations. *Advances in neural information processing systems*, 31, 2018.
- 417 [13] E. Chien, J. Peng, P. Li, and O. Milenkovic. Adaptive universal generalized pagerank graph
418 neural network. In *9th International Conference on Learning Representations, ICLR 2021*,
419 2021.
- 420 [14] F. R. Chung and F. C. Graham. *Spectral graph theory*. Number 92. American Mathematical
421 Soc., 1997.
- 422 [15] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs
423 with fast localized spectral filtering. *Advances in neural information processing systems*, 29,
424 2016.
- 425 [16] J. Eells and J. H. Sampson. Harmonic mappings of riemannian manifolds. *American journal of*
426 *mathematics*, 86(1):109–160, 1964.
- 427 [17] M. Eliasof, E. Haber, and E. Treister. Pde-gcn: Novel architectures for graph neural networks
428 motivated by partial differential equations. *Advances in Neural Information Processing Systems*,
429 34, 2021.
- 430 [18] M. Geiger, L. Petrini, and M. Wyart. Landscape and training regimes in deep learning. *Physics*
431 *Reports*, 924:1–18, 2021.

- 432 [19] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing
433 for quantum chemistry. In *International Conference on Machine Learning*, pages 1263–1272.
434 PMLR, 2017.
- 435 [20] C. Goller and A. Kuchler. Learning task-dependent distributed representations by backprop-
436 agation through structure. In *Proceedings of International Conference on Neural Networks*
437 (*ICNN'96*), volume 1, pages 347–352. IEEE, 1996.
- 438 [21] M. Gori, G. Monfardini, and F. Scarselli. A new model for learning in graph domains. In
439 *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 2,
440 pages 729–734. IEEE, 2005.
- 441 [22] E. Haber and L. Ruthotto. Stable architectures for deep neural networks. *Inverse Problems*, 34,
442 2018.
- 443 [23] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs.
444 *Advances in neural information processing systems*, 30, 2017.
- 445 [24] D. K. Hammond, P. Vandergheynst, and R. Gribonval. The spectral graph wavelet transform:
446 Fundamental theory and fast computation. In *Vertex-Frequency Analysis of Graph Signals*,
447 pages 141–175. Springer, 2019.
- 448 [25] M. He, Z. Wei, H. Xu, et al. Bernnet: Learning arbitrary graph spectral filters via bernstein
449 approximation. *Advances in Neural Information Processing Systems*, 34, 2021.
- 450 [26] R. Kimmel, N. Sochen, and R. Malladi. From high energy physics to low level vision. In
451 *International Conference on Scale-Space Theories in Computer Vision*, pages 236–247. Springer,
452 1997.
- 453 [27] T. N. Kipf and M. Welling. Semi-Supervised Classification with Graph Convolutional Networks.
454 In *Proceedings of the 5th International Conference on Learning Representations, ICLR '17*,
455 2017.
- 456 [28] J. Klicpera, S. Weissenberger, and S. Günnemann. Diffusion improves graph learning. In
457 *Proceedings of the 33rd International Conference on Neural Information Processing Systems*,
458 2019.
- 459 [29] H. Nt and T. Maehara. Revisiting graph neural networks: All we have is low-pass filters. *arXiv*
460 *preprint arXiv:1905.09550*, 2019.
- 461 [30] K. Oono and T. Suzuki. Graph neural networks exponentially lose expressive power for node
462 classification. In *International Conference on Learning Representations*, 2020.
- 463 [31] H. Pei, B. Wei, K. C. Chang, Y. Lei, and B. Yang. Geom-gcn: Geometric graph convolutional
464 networks. In *8th International Conference on Learning Representations, ICLR 2020*, 2020.
- 465 [32] P. Perona and J. Malik. Scale-space and edge detection using anisotropic diffusion. *PAMI*,
466 12(7):629–639, 1990.
- 467 [33] B. Rozemberczki, C. Allen, and R. Sarkar. Multi-scale attributed node embedding. *Journal of*
468 *Complex Networks*, 9(2):cnab014, 2021.
- 469 [34] T. K. Rusch, B. P. Chamberlain, J. Rowbottom, S. Mishra, and M. M. Bronstein. Graph-coupled
470 oscillator networks. In *International Conference on Machine Learning*, 2022.
- 471 [35] M. E. Sander, P. Ablin, M. Blondel, and G. Peyré. Sinkformers: Transformers with doubly
472 stochastic attention. In *International Conference on Artificial Intelligence and Statistics*, pages
473 3515–3530. PMLR, 2022.
- 474 [36] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural
475 network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
- 476 [37] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad. Collective classifica-
477 tion in network data. *AI magazine*, 29(3):93–93, 2008.

- 478 [38] A. Sperduti. Encoding labeled graphs by labeling raam. *Advances in Neural Information*
479 *Processing Systems*, 6, 1993.
- 480 [39] M. Thorpe, T. M. Nguyen, H. Xia, T. Strohmer, A. Bertozzi, S. Osher, and B. Wang. Grand++:
481 Graph neural diffusion with a source term. In *International Conference on Learning Representations*, 2021.
482
- 483 [40] J. Topping, F. Di Giovanni, B. P. Chamberlain, X. Dong, and M. M. Bronstein. Understanding
484 over-squashing and bottlenecks on graphs via curvature. *International Conference on Learning*
485 *Representations*, 2022.
- 486 [41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
487 I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*,
488 30, 2017.
- 489 [42] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph attention
490 networks. In *International Conference on Learning Representations*, 2018.
- 491 [43] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger. Simplifying graph convolutional
492 networks. In *International conference on machine learning*, pages 6861–6871. PMLR, 2019.
- 493 [44] L.-P. Xhonneux, M. Qu, and J. Tang. Continuous graph neural networks. In *International*
494 *Conference on Machine Learning*, pages 10432–10441. PMLR, 2020.
- 495 [45] Y. Yan, M. Hashemi, K. Swersky, Y. Yang, and D. Koutra. Two sides of the same coin:
496 Heterophily and oversmoothing in graph convolutional neural networks. *arXiv preprint*
497 *arXiv:2102.06462*, 2021.
- 498 [46] Z. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec. Gnnexplainer: Generating expla-
499 nations for graph neural networks. *Advances in neural information processing systems*, 32,
500 2019.
- 501 [47] H. Yuan, H. Yu, S. Gui, and S. Ji. Explainability in graph neural networks: A taxonomic survey.
502 *arXiv preprint arXiv:2012.15445*, 2020.
- 503 [48] L. Zhao and L. Akoglu. Pairnorm: Tackling oversmoothing in gnns. *arXiv preprint*
504 *arXiv:1909.12223*, 2019.
- 505 [49] D. Zhou and B. Schölkopf. Regularization on discrete spaces. In *Joint Pattern Recognition*
506 *Symposium*, pages 361–368. Springer, 2005.
- 507 [50] K. Zhou, X. Huang, D. Zha, R. Chen, L. Li, S.-H. Choi, and X. Hu. Dirichlet energy constrained
508 learning for deep graph neural networks. *Advances in Neural Information Processing Systems*,
509 34:21834–21846, 2021.
- 510 [51] J. Zhu, Y. Yan, L. Zhao, M. Heimann, L. Akoglu, and D. Koutra. Beyond homophily in graph
511 neural networks: Current limitations and effective designs. *Advances in Neural Information*
512 *Processing Systems*, 33:7793–7804, 2020.
- 513 [52] O. Shchur, M. Mumme, A. Bojchevski, and S. Günnemann. Pitfalls of graph neural network
514 evaluation. In *NIPS workshop*, 2018.
- 515 [53] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,
516 N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani,
517 S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style,
518 high-performance deep learning library. In *NeurIPS*. 2019.
- 519 [54] M. Fey and J. E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR*
520 *Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- 521 [55] L. Biewald. Experiment tracking with weights and biases, 2020. Software available from
522 wandb.com.

523 Checklist

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

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

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

- 535 1. For all authors...
 - 536 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
537 contributions and scope? **[Yes]**
 - 538 (b) Did you describe the limitations of your work? **[Yes]**, in Section **6**
 - 539 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** in the
540 **Societal impact** paragraph in Section **6**
 - 541 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
542 them? **[Yes]**
- 543 2. If you are including theoretical results...
 - 544 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 545 (b) Did you include complete proofs of all theoretical results? **[Yes]** in Appendix **A**,
546 Appendix **B** and Appendix **C**
- 547 3. If you ran experiments...
 - 548 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
549 imental results (either in the supplemental material or as a URL)? **[Yes]** Code and
550 README in SM, dataloaders in code
 - 551 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
552 were chosen)? **[Yes]** Splits and hyperparameters provided in code zip
 - 553 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
554 ments multiple times)? **[Yes]** Standard deviations are stated in results table
 - 555 (d) Did you include the total amount of compute and the type of resources used (e.g., type
556 of GPUs, internal cluster, or cloud provider)? **[Yes]** in appendix **D**
- 557 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 558 (a) If your work uses existing assets, did you cite the creators? **[Yes]** datasets and standard
559 libraries cited in appendix **D**
 - 560 (b) Did you mention the license of the assets? **[Yes]** industry standard libraries and
561 benchmark datasets were used in accordance with licences
 - 562 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
563 code provided in SM zip
 - 564 (d) Did you discuss whether and how consent was obtained from people whose data you’re
565 using/curating? **[N/A]**
 - 566 (e) Did you discuss whether the data you are using/curating contains personally identifiable
567 information or offensive content? **[Yes]** no personal data is contained within bench-
568 marking datasets
- 569 5. If you used crowdsourcing or conducted research with human subjects...
 - 570 (a) Did you include the full text of instructions given to participants and screenshots, if
571 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]