

---

# How does over-squashing affect the power of GNNs?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Graph Neural Networks (GNNs) are the state-of-the-art model for machine learning  
2 on graph-structured data. The most popular class of GNNs operate by exchanging  
3 information between adjacent nodes, and are known as Message Passing Neural  
4 Networks (MPNNs). While understanding the expressive power of MPNNs is a  
5 key question, existing results typically consider settings with uninformative node  
6 features. In this paper, we provide a rigorous analysis to determine which function  
7 classes of node features can be learned by an MPNN of a given capacity. We do  
8 so by measuring the level of *pairwise interactions* between nodes that MPNNs  
9 allow for. This measure provides a novel quantitative characterization of the so-  
10 called over-squashing effect, which is observed to occur when a large volume  
11 of messages is aggregated into fixed-size vectors. Using our measure, we prove  
12 that, to guarantee sufficient communication between pairs of nodes, the capacity  
13 of the MPNN must be large enough, depending on properties of the input graph  
14 structure, such as commute times. For many relevant scenarios, our analysis results  
15 in impossibility statements in practice, showing that *over-squashing hinders the*  
16 *expressive power of MPNNs*. Our theory also holds for geometric graphs and hence  
17 extends to equivariant MPNNs on point clouds. We validate our analysis through  
18 extensive controlled experiments and ablation studies.

## 19 1 Introduction

20 Graphs describe the relational structure for a large variety of natural and artificial systems, making  
21 learning on graphs imperative in many contexts [48, 20, 51]. Given an underlying graph and features,  
22 defined on its nodes (and edges), as inputs, a Graph Neural Network (GNN) learns parametric  
23 functions from data. Due to the ubiquity of GNNs, characterizing their **expressive power**, i.e., which  
24 class of functions a GNN is able to learn, is a problem of great interest. In this context, most available  
25 results in literature on the universality of GNNs pertain to impractical higher-order tensors [38, 33] or  
26 unique node identifiers that may break the symmetries of the problem [36]. In particular, these results  
27 do not necessarily apply to Message Passing Neural Networks (MPNNs) [27], which have emerged as  
28 the most popular class of GNN models in recent years. Concerning expressivity results for MPNNs,  
29 the most general available characterization is due to [52] and [40], who proved that MPNNs are, at  
30 most, as powerful as the Weisfeiler-Leman graph isomorphism test [50] in distinguishing graphs  
31 *without any features*. This brings us to an important question:

32 Which classes of functions can MPNNs of a given capacity learn, *if node features are specified?*

33 Razin et al. [43] address this question by characterizing the separation rank of MPNNs; however,  
34 their analysis only covers *unconventional* architectures that do not correspond to MPNN models  
35 used in practice. In contrast, Alon & Yahav [2] investigate this question *empirically*, by observing  
36 that MPNNs fail to solve tasks which involve *long-range interactions* among nodes. This limitation  
37 was ascribed to a phenomenon termed as **over-squashing**, which loosely entails messages being

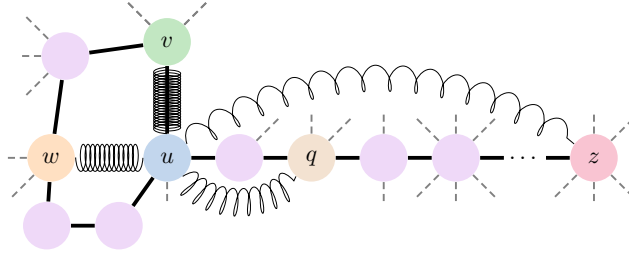


Figure 1: We study the power of MPNNs in terms of the mixing they induce among features and show that this is affected by the model (via norm of the weights and depth) and the graph topology (via commute times). For the given graph, the MPNN learns stronger mixing (tight springs) for nodes  $v, u$  and  $u, w$  since their commute time is small, while nodes  $u, q$  and  $u, z$ , with high commute-time, have weak mixing (loose springs). We characterize over-squashing as the inverse of the mixing induced by an MPNN and hence relate it to its power. In fact, the MPNN might require an impractical depth to solve tasks on the given graph that depend on high-mixing of features assigned to  $u, z$ .

38 ‘squashed’ into fixed-size vectors when the receptive field of a node grows too fast. This effect  
 39 was formalized in [47, 21, 8], who showed that the Jacobian of the nodes features is affected by  
 40 topological properties of the graph, such as curvature and effective resistance. However, all the  
 41 aforementioned papers ignore the specifics of the *task* at hand, i.e., the underlying *function* that the  
 42 MPNN seeks to learn, leading us to the following question:

43 How does *over-squashing affect the expressive power* of MPNNs? Can we *measure* it?

44 **What about geometric graphs?** In many scientific applications, data come as graphs embedded in  
 45 Euclidean space. Since popular architectures resort to the message-passing paradigm [25, 17, 7], the  
 46 expressive power of such models has been rephrased in the language of the WL test, once extended  
 47 to account for the extra geometric information [31]. Nonetheless, the questions raised above are even  
 48 more pressing for these tasks, where the graph is typically derived from a point cloud using a cutoff  
 49 radius, while the features also contain information about the positions in 3D space. In fact, for such  
 50 problems where the features arguably carry more valuable information than the 2D graph structure,  
 51 we argue that proposing new ways to assess the power of message-passing other than (variants of) the  
 52 WL test, is crucial. To this aim, in our paper we study generic message-passing equations with no  
 53 assumptions on the nature of the features, meaning that they may also **include additional positional**  
 54 **information** if the dataset is a point cloud embedded in Euclidean space.

55 **Contributions.** Our main goal is to show how over-squashing can be understood as the *misalignment*  
 56 *between the task and the graph-topology*, ultimately *limiting the classes of functions that MPNNs of*  
 57 *practical size can learn* (see Figure 1). We start by measuring the extent to which an MPNN allows  
 58 pairs of nodes to interact (via **mixing** their features). With this measure as a tool, we characterize  
 59 which functions of node features can be learned by an MPNN and how the model architecture and  
 60 parameters, as well as the topology of the graph, affect the expressive power. More concretely,

- 61 • We introduce a new metric of expressivity based on the Hessian of the function learned by  
 62 an MPNN, which measures the ability of a model to mix features associated with different  
 63 nodes. We then prove upper bounds on the power of MPNNs to mix features (i.e., model  
 64 interactions) according to the novel metric mentioned above. As far as we know, this is the  
 65 first theoretical result stating limitations of MPNNs to learn functions *and their derivatives*.
- 66 • We characterize over-squashing as the reciprocal of the maximal mixing induced by an  
 67 MPNN: *the higher this measure, the smaller the class of functions MPNNs can learn*.
- 68 • We prove that the weights and depth must be sufficiently large – depending on the topology  
 69 – to ensure mixing. For some tasks, *the depth must exceed the highest commute time on the*  
 70 *graph*, resulting in *impossibility* statements. Our results show that MPNNs of practical size,  
 71 fail to learn functions with strong mixing among nodes at high commute time.
- 72 • We illustrate our theoretical results with controlled experiments that verify our analysis, by  
 73 highlighting the impact of the architecture (depth), of the topology (commute time), and of  
 74 the underlying task (the level of mixing required).

75 **2 The Message-Passing paradigm**

76 **Definitions on graphs.** We denote a graph by  $G = (V, E)$ , where  $V$  is the set of  $n$  nodes while  $E$   
 77 are the edges. We assume that  $G$  is *undirected*, *connected* and non-bipartite and define the  $n \times n$   
 78 adjacency matrix  $\mathbf{A}$  as  $A_{vu} = 1$  if  $(v, u) \in E$  and zero otherwise. We let  $\mathbf{D}$  be the diagonal degree  
 79 matrix with  $D_{vv} := d_v$  and use  $d_{\max}$  and  $d_{\min}$  to denote the maximal and minimal degrees. Since  
 80 we are interested in the over-squashing phenomenon, which affects the propagation of information,  
 81 we need to quantify distances on  $G$ . We let  $d_G(v, u)$  be the length of the shortest path connecting  
 82 nodes  $v$  and  $u$  (geodesic distance). While  $d_G$  describes how far two nodes  $u, v$  are in  $G$ , it does not  
 83 account for how many different routes they can use to communicate. In fact, we will see below that  
 84 the over-squashing of nodes  $v, u$  and, more generally, the mixing induced by MPNNs among the  
 85 features associated with  $v, u$ , can be better quantified by their *commute time*  $\tau(v, u)$ , equal to the  
 86 expected number of steps for a random walk to start at  $v$ , reach  $u$ , and then come back to  $v$ .

87 **The MPNN-class.** For most problems, graphs are equipped with features  $\{\mathbf{x}_v\}_{v \in V} \subset \mathbb{R}^d$ , whose  
 88 matrix representation is  $\mathbf{X} \in \mathbb{R}^{n \times d}$ . To study the interactions induced by a GNN among pairs of  
 89 features, we focus on *graph-level* tasks – in Section E of the Appendix, we extend the discussion and  
 90 our main theoretical results to *node-level* tasks. The goal then is to predict a function  $\mathbf{X} \mapsto y_G(\mathbf{X})$ ,  
 91 where we assume that the graph  $G$  is fixed and thus  $y_G : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}$  is a function of the node features.  
 92 MPNNs define a family of parametric functions through iterative local updates of the node features:  
 93 the feature of node  $v$  at layer  $t$  is derived as

$$\mathbf{h}_v^{(t)} = f^{(t)}\left(\mathbf{h}_v^{(t-1)}, g^{(t)}\left(\{\!\!\{\mathbf{h}_u^{(t-1)}, (v, u) \in E\}\!\!\}\right)\right), \quad \mathbf{h}_v^{(0)} = \mathbf{x}_v \quad (1)$$

94 where  $f^{(t)}, g^{(t)}$  are learnable functions and the aggregation function  $g^{(t)}$  is invariant to permutations.  
 95 Specifically, we study a class of MPNNs of the following form,

$$\mathbf{h}_v^{(t)} = \sigma\left(\Omega^{(t)}\mathbf{h}_v^{(t-1)} + \mathbf{W}^{(t)} \sum_u A_{vu} \psi^{(t)}(\mathbf{h}_v^{(t-1)}, \mathbf{h}_u^{(t-1)})\right), \quad \mathbf{h}_v^{(0)} = \mathbf{x}_v, \quad (2)$$

96 where  $\sigma$  acts pointwise,  $\Omega^{(t)}, \mathbf{W}^{(t)} \in \mathbb{R}^{d \times d}$  are weight matrices,  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is any matrix satisfying  
 97  $A_{vu} > 0$  if  $(v, u) \in E$  and zero otherwise –  $\mathbf{A}$  is typically some (normalized) version of the adjacency  
 98 matrix  $\mathbf{A}$  – and  $\psi^{(t)}$  is a learnable message function. The layer-update in (2) includes common  
 99 MPNN-models such as GCN [34], SAGE [28], GIN [52], and GatedGCN [10]. As commented in  
 100 Section 6, this is the most general class of MPNN equations studied thus far in theoretical works on  
 101 over-squashing; unless otherwise stated, all our considerations and analysis apply to MPNNs as in  
 102 (2). For graph-level tasks, a permutation-invariant readout READ is required – usually MAX, MEAN,  
 103 or SUM. We define the graph-level function computed by the MPNN after  $m$  layers to be

$$y_G^{(m)}(\mathbf{X}) = \boldsymbol{\theta}^\top \text{READ}(\{\!\!\{\mathbf{h}_v^{(m)}\}\!\!\}), \quad (3)$$

104 for some learnable  $\boldsymbol{\theta} \in \mathbb{R}^d$ . We restrict to a linear layer since we are interested in the mixing induced  
 105 by the MPNN itself through the topology (and not in readout, independently of the graph-structure).

106 **MPNNs on geometric graphs.** Eq. (2) also describes a class of generic, *equivariant* MPNNs over  
 107 a point cloud embedded in Euclidean space, once the matrix  $\mathbf{A}$  is intended to encode the pairs of  
 108 points that exchange information across each layer. In fact, throughout our analysis, we have no  
 109 restriction on the type of features  $\mathbf{h}_v$ , which can also contain the position of a node in 3D space.  
 110 Accordingly, our theoretical results hold for MPNNs on both 2D and 3D data, since they pertain to  
 111 how the message passing paradigm models pairwise interactions among different points (nodes).

112 **3 On the mixing induced by Message Passing Neural Networks**

113 As one of the main contributions of this paper, we propose a new framework for characterizing the  
 114 expressive power of MPNNs by *estimating the amount of mixing they induce among pairs of node*  
 115 *features*  $\mathbf{x}_v$  and  $\mathbf{x}_u$ . To motivate our definition, fix the underlying graph  $G$ , let  $y_G$  be the ground-truth  
 116 function to be learned, and suppose, for simplicity, that the node features  $\{x_i\}$  are all scalar. If  $y_G$  is  
 117 a smooth function, then we can take the Taylor expansion of  $y_G$  at any point  $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_n)$  and  
 118 obtain a polynomial in the variables  $(x_1, \dots, x_n)$ , up to higher-order corrections. The *mixing* induced

119 by  $y_G$  on the features  $x_v, x_u$  can then be expressed in terms of *mixed* product monomials of the form  
 120  $x_v x_u$ , and the powers thereof. The lowest-degree mixed monomials of this form are multiplied by  
 121 the Hessian (i.e. the second-order derivatives) of  $y_G$ . Accordingly, we can take the entries  $v, u$  of the  
 122 Hessian of  $y_G$  as the simplest **measure of pairwise mixing** induced by  $y_G$  over the nodes  $v, u$ .

123 **Definition 3.1.** For a twice differentiable graph-function  $y_G$  of node features  $\{\mathbf{x}_i\}$ , the **maximal**  
 124 **mixing** induced by  $y_G$  among the features  $\mathbf{x}_v$  and  $\mathbf{x}_u$  associated with nodes  $v, u$  is

$$\text{mix}_{y_G}(v, u) = \max_{\mathbf{x}_i} \max_{1 \leq \alpha, \beta \leq d} \left| \frac{\partial^2 y_G(\mathbf{X})}{\partial x_v^\alpha \partial x_u^\beta} \right|. \quad (4)$$

125 We note that the first maximum is taken over all input features, while the second maximum is taken  
 126 over all entries  $\alpha, \beta$  of the  $d$ -dimensional node features  $\mathbf{x}_v$  and  $\mathbf{x}_u$ ; it is straightforward to adapt the  
 127 results below to alternative definitions based on different norms of the Hessian.

128 **Problem statement.** We study the expressive power of MPNNs in terms of the (maximal) mixing  
 129 they can generate among nodes  $v, u$ . A low value of mixing implies that the MPNN cannot learn  
 130 functions  $y_G$  that require high mixing of the features associated with  $v, u$  and hence it cannot model  
 131 ‘product’-type interactions, as per our explanation above. We investigate how weights and depth on  
 132 the one side, and the graph topology on the other, affect the mixing of an MPNN.

133 **The requirement of smoothness.** In many applications, especially when deploying neural network  
 134 models to solve partial differential equations, the predictions need to be sufficiently regular (smooth),  
 135 which motivates the adoption of smooth activations [29, 9, 23]. Our analysis below follows this  
 136 paradigm and holds for all activations  $\sigma$  that are (at least) twice differentiable.

### 137 3.1 Pairwise mixing induced by MPNNs

138 In this Section, our goal is to derive an upper bound on the maximal mixing induced by MPNNs,  
 139 as defined above, over the features associated with pairs of nodes  $v, u$ . To *motivate the structure of*  
 140 *this bound*, we consider the simple yet illustrative setting of an MPNN as in (2) with scalar features,  
 141 weights  $\omega, w > 0$  and a linear message function of the form  $\psi(x, y) = c_1 x + c_2 y$ , for some learnable  
 142 constants  $c_1, c_2$ . In this case, the layer-update (2) takes the very simple form,

$$h_v^{(t)} = \sigma(w(\mathbf{S}\mathbf{h}^{(t-1)})_v), \quad \mathbf{S} := \frac{\omega}{w}\mathbf{I} + c_1 \text{diag}(\mathbf{A}\mathbf{1}) + c_2 \mathbf{A} \in \mathbb{R}^{n \times n}, \quad (5)$$

143 where  $\mathbf{1} \in \mathbb{R}^n$  is the vector of ones. Hence, the operator  $w\mathbf{S}$  governs the flow of information from  
 144 layer  $t - 1$  to layer  $t$  – once we factor out the derivatives of  $\sigma$  – and the  $k$ -power of this matrix  $(w\mathbf{S})^k$   
 145 determines the propagation of information on the graph over  $k$  layers, i.e. over walks of length  $k$ .

146 A similar argument also works in the general case of (2), once we account for bounds on the non-linear  
 147 activation function  $\sigma$  by  $c_\sigma = \max\{|\sigma'|, |\sigma''|\}$ , and we choose  $\omega, w, c_1, c_2$  satisfying

$$\|\boldsymbol{\Omega}^{(t)}\| \leq \omega, \quad \|\mathbf{W}^{(t)}\| \leq w, \quad \|\nabla_i \psi^{(t)}\| \leq c_i,$$

148 for  $i = 1, 2$ , where  $\nabla_i \psi$  is the Jacobian of  $\psi$  with respect to the  $i$ -th variable, and  $\|\cdot\|$  is the *operator*  
 149 *norm* of a matrix. We note that for trained MPNNs the weights would be finite and bounded, so our  
 150 assumption is mild. For models such as GCN, SAGE or GIN, these constants will suffice in deriving  
 151 the upper bound on the mixing. However, in the general case of non-linear message functions  $\psi$ ,  
 152 which for example includes GatedGCN, we also need to account for the term  $\mathbf{Q}_k$ , defined below,  
 153 which arises when taking second-order derivatives of the MPNN (2): given  $\mathbf{S}$  in (5), we set

$$\begin{aligned} \mathbf{P}_k &:= (\mathbf{S}^{m-k-1})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) (\mathbf{A} \mathbf{S}^{m-k-1}) \\ \mathbf{Q}_k &:= \mathbf{P}_k + \mathbf{P}_k^\top + (\mathbf{S}^{m-k-1})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k (\text{diag}(\mathbf{A}\mathbf{1}) + \mathbf{A})) \mathbf{S}^{m-k-1}. \end{aligned} \quad (6)$$

154 We assume that the Hessian of  $\psi$  is bounded as  $\|\nabla^2 \psi^{(t)}\| \leq c^{(2)}$ . Recall that  $y_G^{(m)}$  is the MPNN-  
 155 prediction (3) and that  $\text{mix}_{y_G^{(m)}}(v, u)$  is its maximal mixing of nodes  $v, u$  as per Definition 3.1.

156 **Theorem 3.2.** Consider an MPNN of depth  $m$  as in (2), where  $\sigma$  and  $\psi^{(t)}$  are  $\mathcal{C}^2$  functions and we  
 157 denote the bounds on their derivatives and on the norm of the weights as above. Let  $\mathbf{S}$  and  $\mathbf{Q}_k$  be

158 defined as in (5) and (6), respectively. If the readout is MAX, MEAN or SUM and  $\theta$  in (3) has unit  
 159 norm, then the mixing  $\text{mix}_{y_G^{(m)}}(v, u)$  induced by the MPNN over the features of nodes  $v, u$  satisfies

$$\text{mix}_{y_G^{(m)}}(v, u) \leq \sum_{k=0}^{m-1} (c_\sigma \mathbf{w})^{2m-k-1} \left( \mathbf{w} (\mathbf{S}^{m-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) \mathbf{S}^{m-k} + c^{(2)} \mathbf{Q}_k \right)_{vu}. \quad (7)$$

160 Theorem 3.2 shows how the mixing induced by an MPNN depends on the model (via regularity of  $\sigma$ ,  
 161 norm of the weights  $\mathbf{w}$ , and depth  $m$ ) and on the graph-topology (via the powers of  $\mathbf{A}$ , which enters  
 162 the definition of  $\mathbf{S}$  in (5)). Our goal now is to expand (7) and relate it to known quantities on the  
 163 graph and show how this can be used to characterize the phenomenon of over-squashing. First, we  
 164 introduce a notion of *capacity* of an MPNN in the spirit of [36].

165 **The capacity of an MPNN.** For simplicity, we assume that  $c_\sigma = 1$ , since this is satisfied by most  
 166 commonly used non-linear activations – it is straightforward to extend the analysis to arbitrary  $c_\sigma$ .

167 **Definition 3.3.** Given an MPNN with  $m$  layers and  $\mathbf{w}$  the maximal operator norm of the weights, we  
 168 say that the pair  $(m, \mathbf{w})$  represents the **capacity** of the MPNN.

169 A larger capacity, by increasing  $m$  or  $\mathbf{w}$ , heuristically implies that the MPNN has more power to  
 170 induce larger mixing among the nodes  $v, u$ . Accordingly, given  $v, u$ , we formulate the problem of  
 171 expressivity as: *what is the capacity required to induce enough mixing  $\text{mix}_{y_G}(v, u)$ ?*

172 **Studying expressivity through derivatives.** In applications to physics and PDEs, we may often  
 173 need the neural-network prediction to also match the derivatives of the ground-truth function [29].  
 174 Theorem 3.2 provides an upper bound on the ability of an MPNN to learn functions with non-trivial  
 175 second-order derivatives among nodes. In particular, (7) shows that the second-order derivatives of  
 176 MPNN predictions as in (2), **cannot** approximate second-order derivatives of graph-functions  $y_G$   
 177 whose associated mixing is larger than the right hand side of (7). Our results are more general than  
 178 the over-squashing problem, and represent, to the best of our knowledge, the first theoretical analysis  
 179 on the limitations of MPNNs to approximate classes of functions and the derivatives thereof.

## 180 4 Over-squashing limits the expressive power of MPNNs

181 Over-squashing was originally described in [2] as the failure of MPNNs to propagate information  
 182 across distant nodes. In fact, [47, 8, 21] showed that over-squashing – quantified by the sensitivity of  
 183 node  $v$  to the input feature at node  $u$  via their Jacobian – is affected by topological properties such as  
 184 curvature and effective resistance. In light of these works, it is evident that over-squashing is related  
 185 to the inability of MPNNs to model interactions among certain nodes, depending on the underlying  
 186 graph topology. Since one can rely on the Taylor expansion of a graph function to measure such  
 187 interactions through the second-order derivatives, i.e. the maximal mixing, we leverage Definition 3.1  
 188 to propose a **novel**, broader, but more accurate, characterization of over-squashing:

189 **Definition 4.1.** Given the prediction  $y_G^{(m)}$  of an MPNN with capacity  $(m, \mathbf{w})$ , we define the **pairwise**  
 190 *over-squashing* of  $v, u$  as

$$\text{OSQ}_{v,u}(m, \mathbf{w}) = \left( \text{mix}_{y_G^{(m)}}(v, u) \right)^{-1}.$$

191 Our notion of over-squashing is a *pairwise* measure over the graph that naturally depends on the  
 192 graph-topology, as well as the capacity of the model. In particular, it captures how over-squashing  
 193 pertains to the ability of the model to mix (induce interactions) between different node features. If  
 194 such maximal mixing is large, then there is no obstruction to exchanging information between the  
 195 given nodes and hence the over-squashing measure would be small; conversely, the over-squashing is  
 196 large precisely when the model struggles to mix features associated with nodes  $v$  and  $u$ .

197 In general though, computing the actual mixing induced by an MPNN may be difficult; we can then  
 198 rely on Theorem 3.2 to derive a proxy for the over-squashing measure that will be used to obtain  
 199 necessary conditions on the capacity of an MPNN to induce a required level of mixing:

200 **Definition 4.2.** Given an MPNN with capacity  $(m, \mathbf{w})$ , we approximate  $\text{OSQ}_{v,u}(m, \mathbf{w})$  by

$$\widetilde{\text{OSQ}}_{v,u}(m, \mathbf{w}) := \left( \sum_{k=0}^{m-1} \mathbf{w}^{2m-k-1} \left( \mathbf{w} (\mathbf{S}^{m-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) \mathbf{S}^{m-k} + c^{(2)} \mathbf{Q}_k \right)_{vu} \right)^{-1}.$$



201 First, note that by Theorem 3.2 we have  $\widetilde{\text{OSQ}}_{v,u}(m, w) \leq \text{OSQ}_{v,u}(m, w)$ . If the network has no  
 202 bandwidth through the weights ( $w = 0$ ), then  $\widetilde{\text{OSQ}}_{v,u}(m, 0) = \infty$ . Besides, the proposed measure  
 203 is infinite (i.e., zero mixing) whenever  $2m < d_G(v, u)$ , which captures the special case of *under-*  
 204 *reaching* for graph-level tasks [5]. We also recall that for simplicity we have taken  $c_\sigma = 1$ , but the  
 205 measure extends to arbitrary non-linear activations  $\sigma$ . Finally, we generalize the characterization of  
 206 OSQ to node-level tasks in Section E of the Appendix.

207 We can rephrase our novel approach to studying expressivity through pairwise mixing, in terms of the  
 208 over-squashing measure and its proxy. By Theorem 3.2 we derive that a **necessary condition for a**  
 209 **smooth MPNN to learn a function  $y_G$  with mixing  $\text{mix}_{y_G}(v, u)$  is**

$$\widetilde{\text{OSQ}}_{v,u}(m, w) < (\text{mix}_{y_G}(v, u))^{-1}. \quad (8)$$

210 An MPNN of given capacity might suffer or not from over-squashing, *depending on the level of*  
 211 *mixing required by the underlying task*. Over-squashing can then be understood as the **misalignment**  
 212 between the task and the underlying topology, as measured by the gap between the maximal mixing  
 213 induced by an MPNN over nodes  $v, u$  and the mixing required by the task.

214 **Strategy.** For a given graph  $G$ , to reduce the value of  $\widetilde{\text{OSQ}}$  and hence satisfy (8), the capacity  $(m, w)$   
 215 must satisfy constraints posed by  $G$  and the choice of  $v, u$ . Since we can increase the capacity by  
 216 taking either larger weights or more layers, we consider these two regimes separately. Below, we  
 217 expand (8) in order to derive minimal requirements on the quantities  $w$  and  $m$  to induce a certain  
 218 level of mixing. For simplicity, we restrict our analysis to the case  $\mathbf{A} = \mathbf{A}_{\text{sym}} := \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$   
 219 and extend the results to  $\mathbf{D}^{-1} \mathbf{A}$  and  $\mathbf{A}$  in Section D of the Appendix.

#### 220 4.1 The case of fixed depth $m$ and variable weights norm $w$

221 To assess the ability of the norm of the weights  $w$  to increase the capacity of an MPNN and hence  
 222 reduce  $\widetilde{\text{OSQ}}$ , we consider the limit case where the depth  $m$  is the minimal required for an MPNN to  
 223 have a non-zero mixing among  $v, u$  (half the shortest-walk distance  $d_G$  between the nodes).

224 **Theorem 4.3.** *Let  $\mathbf{A} = \mathbf{A}_{\text{sym}}$ ,  $r := d_G(v, u)$ ,  $m = \lceil r/2 \rceil$ , and  $q$  be the number of paths of length*  
 225  *$r$  between  $v$  and  $u$ . For an MPNN satisfying Theorem 3.2 with capacity  $(m = \lceil r/2 \rceil, w)$ , we find*  
 226  *$\widetilde{\text{OSQ}}_{v,u}(m, w) \cdot (c_2 w)^r (\mathbf{A}^r)_{vu} \geq 1$ . In particular, if the MPNN generates mixing  $\text{mix}_{y_G}(v, u)$ , then*

$$w \geq \frac{d_{\min}}{c_2} \left( \frac{\text{mix}_{y_G}(v, u)}{q} \right)^{\frac{1}{r}}.$$

227 Theorem 4.3 highlights that if the depth is set as the minimum required for *any* non-zero mixing, then  
 228 the norm of the weights  $w$  has to be large enough depending on the connectivity of  $G$  – recall that  
 229 for models as GCN, we have  $c_2 = 1$ . However, increasing  $w$  is not optimal and may lead to poorer  
 230 generalization capabilities [6, 24]. Besides, controlling the maximal operator norm of the weight  
 231 matrices is not easy, especially from below. We report a few examples in Appendix D.

#### 232 4.2 The case of fixed weights norm $w$ and variable depth $m$

233 We now study the (desirable) setting where  $w$  is bounded, and derive the depth necessary to induce  
 234 mixing of nodes  $v, u$ . Below, we let  $0 = \lambda_0 < \lambda_1 \leq \dots \leq \lambda_{n-1}$  be the eigenvalues of the normalized  
 235 graph Laplacian  $\Delta = \mathbf{I} - \mathbf{A}_{\text{sym}}$ ; we note that  $\lambda_1$  is the spectral gap and  $\lambda_{n-1} < 2$  if  $G$  is not bipartite  
 236 [16]. We also recall that  $d_G$  is the shortest-walk distance and  $\tau$  is the **commute time** (defined in  
 237 Section 2). Finally, if  $d_{\max}$  and  $d_{\min}$  denote the maximal and minimal degrees, respectively, we set  
 238  $\gamma := \sqrt{d_{\max}/d_{\min}}$ .

239 **Theorem 4.4.** *Consider an MPNN satisfying Theorem 3.2, with  $\max\{w, \omega/w + c_1 \gamma + c_2\} \leq 1$ ,*  
 240 *and  $\mathbf{A} = \mathbf{A}_{\text{sym}}$ . If  $\widetilde{\text{OSQ}}_{v,u}(m, w) \cdot (\text{mix}_{y_G}(v, u)) \leq 1$ , i.e. the MPNN generates mixing  $\text{mix}_{y_G}(v, u)$*   
 241 *among the features associated with nodes  $v, u$ , then the number of layers  $m$  satisfies*

$$m \geq \frac{\tau(v, u)}{4c_2} + \frac{|\mathbf{E}|}{\sqrt{d_v d_u}} \left( \frac{\text{mix}_{y_G}(v, u)}{\gamma \mu} - \frac{1}{c_2} \left( \frac{\gamma + |1 - c_2 \lambda^*|^{r-1}}{\lambda_1} + 2 \frac{c^{(2)}}{\mu} \right) \right),$$

242 where  $r = d_G(v, u)$ ,  $\mu = 1 + 2c^{(2)}(1 + \gamma)$  and  $|1 - c_2 \lambda^*| = \max_{0 < \ell \leq n-1} |1 - c_2 \lambda_\ell| < 1$ .

243 Theorem 4.4 provides a *necessary condition* on the depth of an MPNN to induce enough mixing  
244 among nodes  $v, u$ . We see that the MPNN must be sufficiently deep if the task depends on interactions  
245 between nodes at high commute time  $\tau$ . Note that the lower bound on the depth can translate into a  
246 *practical impossibility statement*, since the commute time  $\tau$  can be as large as  $\mathcal{O}(n^3)$  [13].

247 If (i) the graph is such that the commute time between  $v, u$  is large, and (ii) the task depends on  
248 high-mixing of features associated with  $v, u$ , then **over-squashing limits the expressive power** of  
249 MPNNs since the depth has to scale impractically with the graph size  $n$ . In contrast to existing  
250 approaches based on the graph-isomorphism test, our results characterize the expressivity of MPNNs  
251 even when **meaningful (e.g. geometric) features are provided**. In fact, Theorem 4.4 implies that:

252 **Corollary 4.5.** *On a graph with features, MPNNs as in Theorem 4.4 with depth  $m \leq n$ , cannot learn*  
253 *functions that induce high mixing among features of nodes with large commute time.*

254 Since the commute time of two adjacent nodes  $v, u$  equals  $2|E|$  if  $(v, u)$  is a **cut-edge** [1], our result  
255 shows that MPNNs may require  $m = \Omega(|E|)$  to generate enough mixing along a cut-edge, drawing a  
256 connection with [55], where it was shown that most GNNs fail to identify cut-edges on *unattributed*  
257 graphs. In fact, our theoretical results are *more general* than assessing the inability of an MPNN  
258 to solve tasks with long-range interactions, and show how over-squashing can be understood as a  
259 fundamental problem associated with how hard is for an MPNN to exchange information between  
260 nodes that are ‘badly connected’, as per their commute time, whatever this information might be.

261 Theorem 4.4 determines the minimal number of layers required to induce mixing among the *specific*  
262 nodes  $v, u$ . If the depth  $m$  does not satisfy the lower bound in Theorem 4.4, then the mixing induced  
263 among  $v, u$  is **smaller** than  $y_G(v, u)$ . However, increasing the number of layers so to satisfy such  
264 constraint may have a detrimental effect to nodes that have small commute time instead.

## 265 5 Experimental validation of the theoretical results

266 Next, we aim to empirically verify the impact of the graph topology (via commute time  $\tau$ ), the GNN  
267 architecture (depth, norm of weights), and the underlying task (node mixing) on over-squashing, as  
268 predicted by our theory. This, however, requires detailed information about the underlying function to  
269 be learned, which is not readily available in practice. Hence, we perform our empirical test  
270 in a controlled environment, but at the same time, we base our experiments on the *real world*  
271 ZINC chemical dataset [30] and constrain the number of molecular graphs to 12K [22]. Moreover,  
272 we exclude the edge features from this experiment and fix the MPNN size to  $\sim 100\text{K}$  parameters.  
273 However, instead of regressing the constrained solubility based on the molecular input graphs, we  
274 define our own synthetic node features as well as our own target values as follows.

275 Let  $\{G^i\}$  be the set of the 12K ZINC molecular graphs. We set all node features to zero, except for  
276 two, which are set to uniform random numbers  $x_{u^i}^i, x_{v^i}^i$  between 0 and 1 (i.e.,  $x_{u^i}^i, x_{v^i}^i \sim \mathcal{U}(0, 1)$ ) for  
277 all  $i$ . The target is set to  $y^i = \tanh(x_{u^i}^i + x_{v^i}^i)$  for all  $i$ . Hence, the task entails a non-linear mixing  
278 with non-vanishing second derivatives. The two non-zero node features  $x_{u^i}^i, x_{v^i}^i$  are positioned on  
279  $G_i$  according to the commute time  $\tau$ , i.e., for a given  $\alpha \in [0, 1]$ , we choose the nodes  $u^i, v^i$  as the  
280  $\alpha$ -quantile of the  $\tau$ -distribution over  $G_i$ . This grants us a control on the level of commute time of the  
281 underlying mixing ( see Fig. 2). We call this graph dataset the *synthetic ZINC dataset*. We consider  
282 four different MPNN models namely GCN [34], GIN [52], GraphSAGE [28], and GatedGCN [10].  
283 Moreover, we choose the MAX-pooling as the GNN readout, which is supported by Theorem 3.2 and  
284 forces the GNNs to make use of the message-passing in order to learn the mixing.

### 285 5.1 The role of commute time

286 In this task, we empirically analyse the effect of the commute time  $\tau$  of the underlying mixing on  
287 the performance of the MPNNs. To this end, we fix the architecture for all considered MPNNs.  
288 In particular, we set the depth to  $m = \max_i \lceil \text{diam}(G^i)/2 \rceil$ , which happens to be  $m = 11$  for the  
289 considered ZINC 12K graphs, such that the MPNNs are guaranteed *not to underreach*. We further  
290 vary the value of the  $\alpha$ -quantile of the  $\tau$ -distributions over the graphs  $G^i$  between 0 and 1, thus  
291 controlling the level of commute times. According to our theoretical findings in Section 4, the  
292 measure  $\widetilde{\text{OSQ}}_{v,u}$  (Definition 4.2) heavily depends on the commute time  $\tau$  of the underlying mixing  
293 as derived in Theorem 4.4 – we verified this in Appendix Fig. 7. Thus, we would expect the MPNNs

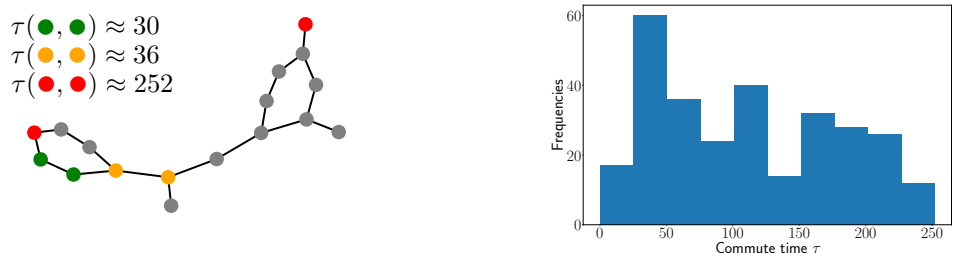


Figure 2: **(Left)** Exemplary molecular graph of the ZINC (12K) dataset with colored nodes corresponding to different values of commute time  $\tau$ . We note that  $\tau$  is a more refined measure than the distance, and in fact beyond long-range nodes (red case),  $\tau$  also captures other topological properties (yellow nodes are adjacent but belong to a *cut-edge*, so their commute-time is  $2|E|$ ). **(Right)** Histogram of commute time  $\tau$  between all pairs of the graph nodes.

294 to perform significantly worse for increasing levels of the commute time. This is indeed confirmed in  
 295 Fig. 3 which shows that the test MAE increases for larger values of  $\alpha$  for all considered MPNNs.

## 296 5.2 The role of depth

297 In this task, we study the effect of the depth on the performance of the MPNNs. To this end, we  
 298 consider a high commute time-regime by setting  $\alpha = 0.8$ . Note that in this case the maximum (over  
 299 all graphs  $G^i$ ) shortest path between two nodes  $u^i, v^i$  is 14. Therefore, a *depth of  $m = 7$  is sufficient*  
 300 *to avoid under-reaching on all graphs*. However, according to the over-squashing measure we provide  
 301 and the conclusions of Theorem 4.4, we expect the MPNNs to be able to induce more mixing among  
 302 nodes  $v, u$ , and hence reduce the error, as we increase the number of layers. This expectation is  
 303 further evidenced in Appendix Fig. 8, where the computed  $\widetilde{\text{OSQ}}$  decreases for increasing number of  
 304 layers. In Fig. 4, we plot the test MAE of all considered MPNNs for increasing number of layers. We  
 305 can indeed see that all considered GNNs benefit from depth, and thus higher capacity (Definition 3.3),  
 306 as GatedGCN obtains the lowest test MAE with 16 layers, as well as GraphSAGE, GIN, and GCN  
 307 with 32 layers. Our theoretical results provide a strong explanation as to why a task **only** depending  
 308 on the mixing of nodes within 14 hops – so that 7 layers would suffice – actually benefits from many  
 309 more layers. Naturally, we cannot increase the depth arbitrarily, as at some point other issues emerge  
 310 which impact the trainability of the MPNNs [44].

311 In Appendix F we also report additional experiments on the role of mixing and how the performance  
 312 of the MPNN models if fully aligned with our theoretical findings.

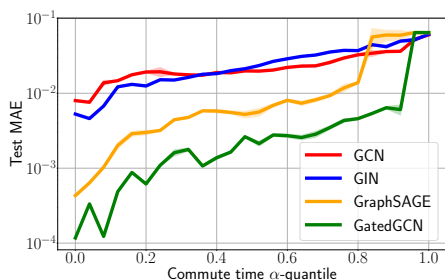


Figure 3: Test MAE (average and standard deviation over several random weight initializations) of GCN, GIN, GraphSAGE, and GatedGCN on synthetic ZINC, where the commute time of the underlying mixing is varied, while the MPNN architecture is fixed, i.e., mixing according to increasing values of the  $\alpha$ -quantile of the  $\tau$ -distribution over the graphs.

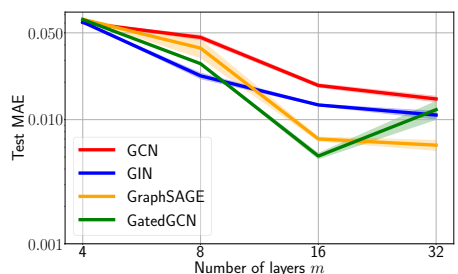


Figure 4: Test MAE (average and standard deviation over several random weight initializations) of GCN, GIN, GraphSAGE, and GatedGCN on synthetic ZINC, where the commute time is fixed to be high (i.e., at the level of the 0.8-quantile), while only the depth of the MPNN is varied between 4 and 32 (all other architectural components are fixed).



313 **6 Discussion**

314 **Related Work: expressive power of MPNNs.** The MPNN class in (2) is as **powerful** as the 1-WL  
 315 test [50] in distinguishing *unattributed* graphs [52, 40]. In fact, MPNNs typically struggle to compute  
 316 graph properties on feature-less graphs [19, 15, 45, 36]. The expressivity of GNNs has also been  
 317 studied from the point of view of logical and tensor languages [5, 4, 26]. Nonetheless, far less is  
 318 known about which functions of node features MPNNs can learn and *the capacity required to do so*.  
 319 Razin et al. [43] recently studied the separation rank of a specific MPNN class. While this approach  
 320 is a strong inspiration for our work, the results in [43] only apply to a family of MPNNs which does  
 321 not include models used in practice. *Our results instead hold in the full generality of (2) and provide*  
 322 *a novel approach for investigating the expressivity of MPNNs through the mixing they are able to*  
 323 *generate among features. To the best of our knowledge, this is the first work formally analysing the*  
 324 *limitations on the expressive power of MPNNs to learn functions and their second-order derivatives.*

325 **Differences between mixing and the WL test.** Throughout our analysis we had **no** assumption on  
 326 the nature of the features, that can in fact be structural or positional – meaning that the MPNNs we  
 327 have considered above, may also be more powerful than the 1-WL test. Our derivations do not rely  
 328 on the ability to distinguish different node representations, but rather on the ability of the MPNN to  
 329 mix information associated with different nodes. This novel alternative paradigm may help design  
 330 GNNs that are more powerful at mixing than MPNNs, and may further shed light on how and when  
 331 frameworks such as Transformers can solve the underlying task better than conventional MPNNs.

332 **Differences between our results and existing works on over-squashing.** The problem of **over-**  
 333 **squashing** was introduced in [2] and studied through sensitivity analysis in [47]. This approach  
 334 was generalized in [8, 21] who proved that the Jacobian of node features is likely to be small if  
 335 the nodes have high commute time (effective resistance). We discuss more in detail the novelty  
 336 of this work when compared to [47, 21]. (i) In [47, 21] there is no analysis on which functions  
 337 MPNNs cannot learn as a consequence of over-squashing, nor a formal measure of over-squashing.  
 338 Besides, the Jacobian of node features may not be suited for studying over-squashing for graph-level  
 339 tasks. Note that our theory also holds for node-level tasks – see Section E of the Appendix. (ii) The  
 340 analysis in [47] does not address over-squashing among nodes at distance larger than 2 and does not  
 341 provide insights on the capacity required to learn certain tasks. (iii) Finally, [21] does not account for  
 342 MPNNs such as GatedGCN (while ours does), and the connection to commute time is only carried  
 343 out under simplifying assumptions on the nonlinear activation. *We have extended these ideas to*  
 344 *connect over-squashing and expressive power by studying higher-order derivatives of the MPNN and*  
 345 *relating them to the capacity of the model and the underlying graph-topology.*

346 **The measures OSQ and  $\widetilde{\text{OSQ}}$ .** Definition 4.1 considers pairs of nodes and second-order derivatives;  
 347 this could be generalized to a hierarchy accounting for higher-order interactions of nodes. Besides, if,  
 348 depending on the problem, one has access to better estimates on the mixing induced by an MPNN  
 349 than (7), then one can extend our approach and get a finer approximation of OSQ.

350 **Beyond sum-aggregations.** Our results apply to MPNNs as in (2), where  $\mathbf{A}$  is constant, and  
 351 do not include attention-based MPNNs [49, 11] or Graph-Transformers [35, 39, 54, 42] which  
 352 further depend on features via normalization. Extending the analysis to these models is only more  
 353 technically involved. More generally, one could replace the aggregation  $\sum_u A_{vu}$  in (2) with a smooth,  
 354 permutation invariant operator  $\oplus$  [12, 41]. *Our formalism will then prove useful to assess if different*  
 355 *aggregations are more expressive in terms of the mixing (interactions) they are able to generate.*

356 **Graph rewiring.** Another way of going beyond (2) to find MPNNs with lower OSQ is to replace  
 357  $\mathbf{A}$  with a different matrix  $\mathbf{A}'$ , (partly) independent of the connectivity of the input graph, obtained  
 358 from some ‘rewiring’ procedure. Theorem 4.4 validates why recent graph-rewiring methods such as  
 359 [3, 32, 18, 8] manage to alleviate over-squashing: by adding edges that decrease the overall effective  
 360 resistance (commute time) of the graph, these methods reduce the measure OSQ. More generally,  
 361 Definition 4.2 allows one to measure whether a given rewiring is beneficial in terms of over-squashing  
 362 (and hence of the mixing generated) and to what extent. In fact, it follows from Theorem 4.4 that  
 363 methods like [18, 46] are in this sense **optimal**, since they propagate information over *expander*  
 364 graphs, which are sparse and have commute time scaling linearly with the number of edges. Finally,  
 365 our results suggest that for data given by point clouds, *the choice of a computational graph over*  
 366 *which message passing can operate, should also account for the commute time associated with it,*  
 367 *given that the latter represents the correct metric to assess over-squashing.*

## 368 References

- 369 [1] Romas Aleliunas, Richard M Karp, Richard J Lipton, László Lovász, and Charles Rackoff.  
370 Random walks, universal traversal sequences, and the complexity of maze problems. In *20th*  
371 *Annual Symposium on Foundations of Computer Science (sfcs 1979)*, pp. 218–223. IEEE  
372 Computer Society, 1979.
- 373 [2] Uri Alon and Eran Yahav. On the bottleneck of graph neural networks and its practical  
374 implications. In *International Conference on Learning Representations*, 2021.
- 375 [3] Adrián Arnaiz-Rodríguez, Ahmed Begga, Francisco Escolano, and Nuria Oliver. DiffWire:  
376 Inductive Graph Rewiring via the Lovász Bound. In *The First Learning on Graphs Conference*,  
377 2022. URL <https://openreview.net/pdf?id=IXvfIex0mX6f>.
- 378 [4] Waiss Azizian and Marc Lelarge. Expressive power of invariant and equivariant graph neural  
379 networks. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=lxHgXYN4bwl>.
- 381 [5] Pablo Barceló, Egor V Kostylev, Mikael Monet, Jorge Pérez, Juan Reutter, and Juan Pablo Silva.  
382 The logical expressiveness of graph neural networks. In *International Conference on Learning*  
383 *Representations*, 2019.
- 384 [6] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds  
385 for neural networks. *Advances in neural information processing systems*, 30, 2017.
- 386 [7] Simon Batzner, Albert Musaelian, Lixin Sun, Mario Geiger, Jonathan P Mailoa, Mordechai  
387 Kornbluth, Nicola Molinari, Tess E Smidt, and Boris Kozinsky. E (3)-equivariant graph neural  
388 networks for data-efficient and accurate interatomic potentials. *Nature communications*, 13(1):  
389 2453, 2022.
- 390 [8] Mitchell Black, Amir Nayyeri, Zhengchao Wan, and Yusu Wang. Understanding oversquashing  
391 in gnns through the lens of effective resistance. *arXiv preprint arXiv:2302.06835*, 2023.
- 392 [9] Johannes Brandstetter, Daniel E Worrall, and Max Welling. Message passing neural pde solvers.  
393 In *International Conference on Learning Representations*, 2021.
- 394 [10] Xavier Bresson and Thomas Laurent. Residual gated graph convnets. *arXiv preprint*  
395 *arXiv:1711.07553*, 2017.
- 396 [11] Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks? In  
397 *International Conference on Learning Representations*, 2022.
- 398 [12] Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning:  
399 Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- 400 [13] Ashok K Chandra, Prabhakar Raghavan, Walter L Ruzzo, Roman Smolensky, and Prason  
401 Tiwari. The electrical resistance of a graph captures its commute and cover times. *computational*  
402 *complexity*, 6(4):312–340, 1996.
- 403 [14] Rongqin Chen, Shenghui Zhang, Ye Li, et al. Redundancy-free message passing for graph  
404 neural networks. *Advances in Neural Information Processing Systems*, 35:4316–4327, 2022.
- 405 [15] Zhengdao Chen, Lei Chen, Soledad Villar, and Joan Bruna. Can graph neural networks count  
406 substructures? *Advances in neural information processing systems*, 33:10383–10395, 2020.
- 407 [16] Fan RK Chung and Fan Chung Graham. *Spectral graph theory*. American Mathematical Soc.,  
408 1997.
- 409 [17] Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F  
410 Milles, Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep  
411 learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- 412 [18] Andreea Deac, Marc Lackenby, and Petar Veličković. Expander graph propagation. In *The First*  
413 *Learning on Graphs Conference*, 2022.

- 414 [19] Nima Dehmamy, Albert-László Barabási, and Rose Yu. Understanding the representation  
415 power of graph neural networks in learning graph topology. *Advances in Neural Information*  
416 *Processing Systems*, 32, 2019.
- 417 [20] Gage DeZoort, Peter W Battaglia, Catherine Biscarat, and Jean-Roch Vlimant. Graph neural  
418 networks at the large hadron collider. *Nature Reviews Physics*, pp. 1–23, 2023.
- 419 [21] Francesco Di Giovanni, Lorenzo Giusti, Federico Barbero, Giulia Luise, Pietro Lio, and Michael  
420 Bronstein. On over-squashing in message passing neural networks: The impact of width, depth,  
421 and topology. In *International Conference on Machine Learning*, 2023.
- 422 [22] Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas Laurent, Yoshua Bengio, and Xavier  
423 Bresson. Benchmarking graph neural networks. *arXiv:2003.00982*, 2020.
- 424 [23] Léonard Equer, T Konstantin Rusch, and Siddhartha Mishra. Multi-scale message passing  
425 neural pde solvers. *arXiv preprint arXiv:2302.03580*, 2023.
- 426 [24] Vikas Garg, Stefanie Jegelka, and Tommi Jaakkola. Generalization and representational limits  
427 of graph neural networks. In *International Conference on Machine Learning*, pp. 3419–3430.  
428 PMLR, 2020.
- 429 [25] Johannes Gasteiger, Florian Becker, and Stephan Günnemann. Gemnet: Universal directional  
430 graph neural networks for molecules. *Advances in Neural Information Processing Systems*, 34:  
431 6790–6802, 2021.
- 432 [26] Floris Geerts and Juan L Reutter. Expressiveness and approximation properties of graph  
433 neural networks. In *International Conference on Learning Representations*, 2022. URL  
434 <https://openreview.net/forum?id=wIzUeM3TAU>.
- 435 [27] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural  
436 message passing for quantum chemistry. In *International Conference on Machine Learning*, pp.  
437 1263–1272. PMLR, 2017.
- 438 [28] W. L. Hamilton, R. Ying, and J. Leskovec. Inductive representation learning on large graphs. In  
439 *Advances in Neural Information Processing Systems*, pp. 1025–1035, 2017.
- 440 [29] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Universal approximation of an  
441 unknown mapping and its derivatives using multilayer feedforward networks. *Neural networks*,  
442 3(5):551–560, 1990.
- 443 [30] John J Irwin, Teague Sterling, Michael M Mysinger, Erin S Bolstad, and Ryan G Coleman. Zinc:  
444 a free tool to discover chemistry for biology. *Journal of chemical information and modeling*, 52  
445 (7):1757–1768, 2012.
- 446 [31] Chaitanya K Joshi, Cristian Bodnar, Simon V Mathis, Taco Cohen, and Pietro Lio. On the  
447 expressive power of geometric graph neural networks. *arXiv preprint arXiv:2301.09308*, 2023.
- 448 [32] Kedar Karhadkar, Pradeep Kr. Banerjee, and Guido Montufar. FoSR: First-order spectral  
449 rewiring for addressing oversquashing in GNNs. In *The Eleventh International Confer-*  
450 *ence on Learning Representations*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=3YjQfCLdrzz)  
451 [3YjQfCLdrzz](https://openreview.net/forum?id=3YjQfCLdrzz).
- 452 [33] Nicolas Keriven and Gabriel Peyré. Universal invariant and equivariant graph neural networks.  
453 *Advances in Neural Information Processing Systems*, 32, 2019.
- 454 [34] Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional  
455 Networks. In *International Conference on Learning Representations*, 2017.
- 456 [35] Devin Kreuzer, Dominique Beaini, Will Hamilton, Vincent Létourneau, and Prudencio Tossou.  
457 Rethinking graph transformers with spectral attention. In *Advances in Neural Information*  
458 *Processing Systems*, volume 34, pp. 21618–21629, 2021.
- 459 [36] Andreas Loukas. What graph neural networks cannot learn: depth vs width. In *International*  
460 *Conference on Learning Representations*, 2020. URL [https://openreview.net/forum?](https://openreview.net/forum?id=B112bp4YwS)  
461 [id=B112bp4YwS](https://openreview.net/forum?id=B112bp4YwS).

- 462 [37] László Lovász. Random walks on graphs. *Combinatorics, Paul erdos is eighty*, 2(1-46):4, 1993.
- 463 [38] Haggai Maron, Ethan Fetaya, Nimrod Segol, and Yaron Lipman. On the universality of invariant  
464 networks. In *International conference on machine learning*, pp. 4363–4371. PMLR, 2019.
- 465 [39] Grégoire Mialon, Dexiong Chen, Margot Selosse, and Julien Mairal. Graphit: Encoding graph  
466 structure in transformers. *CoRR*, abs/2106.05667, 2021.
- 467 [40] Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen,  
468 Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural  
469 networks. In *AAAI Conference on Artificial Intelligence*, pp. 4602–4609. AAAI Press, 2019.
- 470 [41] Euan Ong and Petar Veličković. Learnable commutative monoids for graph neural networks.  
471 *arXiv preprint arXiv:2212.08541*, 2022.
- 472 [42] Ladislav Rampasek, Mikhail Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and  
473 Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. In *Advances in*  
474 *Neural Information Processing Systems*, 2022.
- 475 [43] Noam Razin, Tom Verbin, and Nadav Cohen. On the ability of graph neural networks to model  
476 interactions between vertices. *arXiv preprint arXiv:2211.16494*, 2022.
- 477 [44] T. Konstantin Rusch, Ben Chamberlain, James Rowbottom, Siddhartha Mishra, and Michael  
478 Bronstein. Graph-coupled oscillator networks. In *Proceedings of the 39th International*  
479 *Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,  
480 pp. 18888–18909. PMLR, 2022.
- 481 [45] Ryoma Sato, Makoto Yamada, and Hisashi Kashima. Approximation ratios of graph neural  
482 networks for combinatorial problems. *Advances in Neural Information Processing Systems*, 32,  
483 2019.
- 484 [46] Hamed Shirzad, Ameya Velingker, Balaji Venkatachalam, Danica J Sutherland, and Ali Kemal  
485 Sinop. Exphormer: Sparse transformers for graphs. *arXiv preprint arXiv:2303.06147*, 2023.
- 486 [47] Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, and  
487 Michael M Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature.  
488 In *International Conference on Learning Representations*, 2022.
- 489 [48] Petar Veličković. Everything is connected: Graph neural networks. *Current Opinion in*  
490 *Structural Biology*, 79:102538, 2023.
- 491 [49] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua  
492 Bengio. Graph attention networks. In *International Conference on Learning Representations*,  
493 2018.
- 494 [50] Boris Weisfeiler and Andrei Leman. The reduction of a graph to canonical form and the algebra  
495 which appears therein. *nti, Series*, 2(9):12–16, 1968.
- 496 [51] Geordie Williamson. Is deep learning a useful tool for the pure mathematician? *arXiv preprint*  
497 *arXiv:2304.12602*, 2023.
- 498 [52] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? In  
499 *International Conference on Learning Representations*, 2019.
- 500 [53] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and  
501 Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In  
502 *International Conference on Machine Learning*, pp. 5453–5462. PMLR, 2018.
- 503 [54] Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen,  
504 and Tie-Yan Liu. Do transformers really perform badly for graph representation? In *Advances*  
505 *in Neural Information Processing Systems*, volume 34, pp. 28877–28888, 2021.
- 506 [55] Bohang Zhang, Shengjie Luo, Liwei Wang, and Di He. Rethinking the expressive power of gnns  
507 via graph biconnectivity. In *The Eleventh International Conference on Learning Representations*,  
508 2023.

509 **A Outline of the appendix**

510 We provide an overview of the appendix. Since in the appendix we report additional theoretical  
 511 results and considerations, we first point out to the most relevant content: the proofs of the main  
 512 results, the extension of our discussion and analysis to node-level tasks, and the additional ablation  
 513 studies.

514 **Where to find proofs of the main results.** We prove Theorem 3.2 in Section C.1, we prove Theorem  
 515 4.3 in Section D.1, and finally we prove Theorem 4.4 in Section D.3.

516 **Where to find the extension to node-level tasks.** Concerning the case of node-level tasks, we present  
 517 a thorough discussion on the matter in Section E, where we extend the definition of the over-squashing  
 518 measure and generalize Theorem 3.2 and Theorem 4.4 to node-level predictions of the MPNN class  
 519 in (2).

520 **Where to find additional ablation studies.** In Section F we have conducted further experiments  
 521 on the profile of the over-squashing measure  $\widehat{\text{OSQ}}$  across different MPNN models as well as on the  
 522 training mean average error, to further validate our claims on over-squashing hindering the expressive  
 523 power of MPNNs.

524 Next, we summarize the contents of the Appendix more in detail below.

- 525 • In order to be self-consistent, in Section B we review important notions pertaining to the  
 526 spectrum of the graph-Laplacian and known properties of random walks on graphs, that will  
 527 be then be used in our proofs.
- 528 • In Section C we prove the main theorem on the maximal mixing induced by MPNNs  
 529 (Theorem 3.2). In particular, we also derive additional results on the mixing generated  
 530 at a specific node, which will turn out useful when extending the characterization of the  
 531 over-squashing measure  $\widehat{\text{OSQ}}$  for node-level tasks.
- 532 • In Section D we prove the main results of Section 4, mainly Theorem 4.3 and Theorem 4.4.  
 533 Further, we also derive an explicit (sharper) characterization of the depth required to induce  
 534 enough mixing among nodes, in terms of the pseudo-inverse of the graph-Laplacian. Finally,  
 535 in Section D.4 we extend the results to the case of the unnormalized adjacency matrix and  
 536 discuss relative over-squashing measures.
- 537 • In Section E we generalize the over-squashing measure for node-level tasks, commenting on  
 538 the differences between our approach and existing works (mainly [47, 8, 21]). In particular,  
 539 we show that the same conclusions of Theorem 4.4 hold for node-level predictions too.
- 540 • Finally, in Section F we report additional details on our experimental setup and further  
 541 ablation studies concerning the over-squashing measure  $\widehat{\text{OSQ}}$ .

542 **B Summary of spectral properties on graphs**

543 **Basic notions of spectral theory on graphs.**

544 Throughout the appendix, we let  $\Delta$  be the normalized graph Laplacian defined by  $\Delta = \mathbf{I} -$   
 545  $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ . It is known [16] that the graph Laplacian is a symmetrically, positive semi-definite  
 546 matrix whose spectral decomposition takes the form

$$\Delta = \sum_{\ell=0}^{n-1} \lambda_{\ell} \phi_{\ell} \phi_{\ell}^{\top}, \tag{9}$$

547 where  $\{\phi_{\ell}\}$  is an orthonormal basis in  $\mathbb{R}^n$  and  $0 = \lambda_0 < \lambda_1 < \dots < \lambda_{n-1}$  – recall that since we  
 548 assume  $G$  to be connected, the zero eigenvalue has multiplicity one, i.e.  $\lambda_1 > 0$ . We also note that  
 549 we typically write  $\phi_{\ell}(v)$  for the value of  $\phi_{\ell}$  at  $v \in V$ , and that the kernel of  $\Delta$  is spanned by  $\phi_0$   
 550 with  $\phi_0(v) = \sqrt{d_v/2|E|}$ . the results would extend to the bipartite case as usual when doing spectral  
 551 analysis one too if graphs, we exclude the edge case of the bipartite graph to make sure that the  
 552 largest eigenvalue of the graph Laplacian satisfies  $\lambda_{n-1} < 2$  – yet all results hold for the bipartite  
 553 case too provided we take  $\|\nabla_2 \psi\| < 1$ . Finally, we let  $\Delta^{\dagger}$  denote the pseudo-inverse of the graph



554 Laplacian, which can be written as

$$\Delta^\dagger = \sum_{\ell=1}^{n-1} \frac{1}{\lambda_\ell} \phi_\ell \phi_\ell^\top, \quad (10)$$

555 and we emphasize that the sum starts from  $\ell = 1$  since we need to ignore the kernel of  $\Delta$  spanned by  
556 the orthonormal vector  $\phi_0$ .

557 **Basic properties of Random Walks on graphs.** A simple Random Walk (RW) on  $G$  is a Markov  
558 chain supported on the nodes  $V$  with transition matrix defined by  $P(v, u) = d_v^{-1}$ . While a RW can  
559 be studied through different properties, the one we are interested in is the *commute time*  $\tau$ , which  
560 represents the expected number of steps for a RW starting at  $v$ , to visit  $u$  and then come back to  $v$ .  
561 The commute time is a *distance* on the graph and captures the diffusion properties associated with the  
562 underlying topology. In fact, while nodes that are distant often have larger commute time, the latter is  
563 more expressive than the shortest-walk graph-distance, since it also accounts (for example) for the  
564 number of paths connecting two given nodes. Thanks to [37], we can write down the commute time  
565 among two nodes using the spectral representation of the graph Laplacian in (9):

$$\tau(v, u) = 2|E| \sum_{\ell=1}^{n-1} \frac{1}{\lambda_\ell} \left( \frac{\phi_\ell(v)}{\sqrt{d_v}} - \frac{\phi_\ell(u)}{\sqrt{d_u}} \right)^2. \quad (11)$$

### 566 C Proofs and additional details of Section 3

567 The goal of this section amounts to proving Theorem 3.2. To work towards this result, we first derive  
568 bounds on the Jacobian and Hessian of a single node feature after  $m$  layers before the readout READ  
569 operation. We emphasize that our analysis below is novel, when compared to previous works of  
570 [8, 21], on many accounts. First, [8, 21] do not consider higher (second) order derivatives, limiting  
571 their discussion to the case of first order derivatives, which are not suited to capture notions of mixing  
572 among features – we will expand on this topic in Section E. Second, even for the case of first-order  
573 derivatives, our result below is more general since it holds for all MPNNs as in (2), which includes  
574 (i) message-functions  $\psi$  that also depend on the input features (as for GatedGCN), and (ii) choices  
575 of message-passing matrices  $\mathbf{A}$  that could be weighted and (or) asymmetric. Third, the analysis in  
576 [8, 21] does not account for the role of the readout map and hence fails to study the expressive power  
577 of graph-level prediction of MPNNs as measured by the mixing they generate among nodes.

578 **Conventions and notations for the proofs.** First, we recall that  $\mathbf{h}_v^{(0)} = \mathbf{x}_v \in \mathbb{R}^d$  is the input feature  
579 at node  $v$ . Below, we write  $h_v^{(t),\alpha}$  for the  $\alpha$ -th entry of the feature  $\mathbf{h}_v^{(t)}$ . To simplify the notations, we  
580 rewrite the layer-update in (2) using coordinates as

$$h_v^{(t),\alpha} = \sigma(\tilde{h}_v^{(t-1),\alpha}), \quad 1 \leq \alpha \leq d, \quad (12)$$

581 where  $\tilde{h}_v^{(t-1),\alpha}$  is the entry  $\alpha$  of the pre-activated feature of node  $v$  at layer  $t$ . We also let  $\partial_{1,p}\psi^{(t),r}$   
582 and  $\partial_{2,p}\psi^{(t),r}$  be the  $p$ -th derivative of  $(\psi^{(t)}(\cdot, x))_r$  and of  $(\psi^{(t)}(x, \cdot))_r$ , respectively. To avoid  
583 cumbersome notations, we usually omit to write the arguments of the derivatives of the message-  
584 functions  $\psi$ . Similarly, we let  $\nabla_1\psi$  ( $\nabla_2\psi$ ) be the  $d \times d$  Jacobian matrix of  $\psi$  with respect to the first  
585 (second) variable. Finally, given nodes  $i, v, u \in V$  we introduce the following terms:

$$\nabla_u \mathbf{h}_v^{(m)} := \frac{\partial \mathbf{h}_v^{(m)}}{\partial \mathbf{x}_u} \in \mathbb{R}^{d \times d}, \quad \nabla_{uv}^2 \mathbf{h}_i^{(m)} := \frac{\partial^2 \mathbf{h}_i^{(m)}}{\partial \mathbf{x}_u \partial \mathbf{x}_v} \in \mathbb{R}^{d \times (d \times d)}.$$

586 First, we derive an upper bound on the first-order derivatives of the node-features. This will provide  
587 useful to derive the more general second-order estimate of the MPNN-prediction. We highlight that  
588 the result below extends the analysis in [21] to MPNNs with arbitrary (i.e. non-linear) message  
589 functions  $\psi$ , such as GatedGCN [10].

590 **Theorem C.1.** *Given MPNNs as in (2), let  $\sigma$  and  $\psi^{(t)}$  be  $\mathcal{C}^1$  functions and assume  $|\sigma'| \leq c_\sigma$ ,  
591  $\|\mathbf{\Omega}^{(t)}\| \leq \omega$ ,  $\|\mathbf{W}^{(t)}\| \leq w$ ,  $\|\nabla_1\psi^{(t)}\| \leq c_1$ , and  $\|\nabla_2\psi^{(t)}\| \leq c_2$ . Let  $\mathbf{S} \in \mathbb{R}^{n \times n}$  be*

$$\mathbf{S} := \frac{\omega}{w} \mathbf{I} + c_1 \text{diag}(\mathbf{A}\mathbf{1}) + c_2 \mathbf{A}.$$

592 *Given nodes  $v, u \in V$  and  $m$  the number of layers, the following holds:*

$$\|\nabla_u \mathbf{h}_v^{(m)}\| \leq (c_\sigma w)^m (\mathbf{S}^m)_{vu}. \quad (13)$$

593 *Proof.* Recall that the dimension of the features is taken to be  $d$  for any layer  $1 \leq t \leq m$ . We proceed  
 594 by induction. If  $m = 1$  and we fix entries  $1 \leq \alpha, \beta \leq d$ , then using the shorthand in (12), we obtain

$$\begin{aligned} (\nabla_u \mathbf{h}_v^{(1)})_{\alpha\beta} &= \sigma'(\tilde{h}_v^{(0),\alpha}) \left( \Omega_{\alpha\beta}^{(1)} \delta_{vu} + \sum_r W_{\alpha r}^{(1)} \sum_j \mathbf{A}_{vj} \left( \partial_{1,\beta} \psi^{(1),r} \delta_{vu} + \partial_{2,\beta} \psi^{(1),r} \delta_{ju} \right) \right) \\ &= \left( \text{diag}(\sigma'(\tilde{\mathbf{h}}_v^{(0)})) \left( \Omega^{(1)} \delta_{vu} + \mathbf{W}^{(1)} \left( \sum_j \mathbf{A}_{vj} \delta_{vu} \nabla_1 \psi^{(1)} + \mathbf{A}_{vu} \nabla_2 \psi^{(1)} \right) \right) \right)_{\alpha\beta}. \end{aligned}$$

595 Therefore, we can bound the (spectral) norm of the Jacobian on the left hand side by

$$\begin{aligned} \|\nabla_u \mathbf{h}_v^{(1)}\| &\leq \|\text{diag}(\sigma'(\tilde{\mathbf{h}}_v^{(0)}))\| \left( \|\Omega^{(1)}\| \delta_{vu} + \|\mathbf{W}^{(1)}\| (c_1 \sum_j \mathbf{A}_{vj} \delta_{vu} + c_2 \mathbf{A}_{vu}) \right) \\ &\leq c_\sigma (\omega \delta_{vu} + w (c_1 \sum_j \mathbf{A}_{vj} \delta_{vu} + c_2 \mathbf{A}_{vu})) = c_\sigma w \mathbf{S}_{vu}, \end{aligned}$$

596 which proves the estimate on the Jacobian for the case of  $m = 1$ . We now take the induction step,  
 597 and follow the same argument above to write the node Jacobian after  $m$  layers as

$$\begin{aligned} (\nabla_u \mathbf{h}_v^{(m)})_{\alpha\beta} &= \sigma'(\tilde{h}_v^{(m-1),\alpha}) \left( \sum_r \Omega_{\alpha r}^{(m)} (\nabla_u \mathbf{h}_v^{(m-1)})_{r\beta} \right) \\ &\quad + \sigma'(\tilde{h}_v^{(m-1),\alpha}) \left( W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{vj} \sum_p \left( \partial_{1,p} \psi^{(m),r} (\nabla_u \mathbf{h}_v^{(m-1)})_{p\beta} + \partial_{2,p} \psi^{(m),r} (\nabla_u \mathbf{h}_j^{(m-1)})_{p\beta} \right) \right) \\ &= \left( \text{diag}(\sigma'(\tilde{\mathbf{h}}_v^{(m-1)})) \Omega^{(m)} \nabla_u \mathbf{h}_v^{(m-1)} \right)_{\alpha\beta} \\ &\quad + \left( \text{diag}(\sigma'(\tilde{\mathbf{h}}_v^{(m-1)})) \mathbf{W}^{(m)} \left( \sum_j \mathbf{A}_{vj} \nabla_1 \psi^{(m)} \nabla_u \mathbf{h}_v^{(m-1)} + \mathbf{A}_{vj} \nabla_2 \psi^{(m)} \nabla_u \mathbf{h}_j^{(m-1)} \right) \right)_{\alpha\beta}. \end{aligned}$$

598 Therefore, we can use the induction step to bound the Jacobian as

$$\begin{aligned} \|\nabla_u \mathbf{h}_v^{(m)}\| &\leq c_\sigma \omega (c_\sigma w)^{m-1} (\mathbf{S}^{m-1})_{vu} + (c_\sigma w)^m \left( \sum_j \mathbf{A}_{vj} c_1 (\mathbf{S}^{m-1})_{vu} + \mathbf{A}_{vj} c_2 (\mathbf{S}^{m-1})_{ju} \right) \\ &= (c_\sigma w)^m \left( \left( \frac{\omega}{w} \mathbf{I} + c_1 \text{diag}(\mathbf{A}\mathbf{1}) + c_2 \mathbf{A} \right) (\mathbf{S}^{m-1}) \right)_{vu} = (c_\sigma w)^m (\mathbf{S}^m)_{vu}, \end{aligned}$$

599 which completes the proof for the first-order bounds.  $\square$

600 Before we move to the second-order estimates, we introduce some additional preliminary notations.  
 601 Given nodes  $i, v, u$ , a matrix  $\mathbf{S} \in \mathbb{R}^{n \times n}$  – which will always be chosen as per (5) – and an integer  $\ell$ ,  
 602 we write

$$\mathbf{P}_{i(vu)}^{(\ell)} := (\mathbf{S}^\ell)_{iv} (\mathbf{A}\mathbf{S}^\ell)_{iu} + (\mathbf{S}^\ell)_{iu} (\mathbf{A}\mathbf{S}^\ell)_{iv} + \sum_j (\mathbf{S}^\ell)_{jv} (\text{diag}(\mathbf{A}\mathbf{1}) + \mathbf{A})_{ij} (\mathbf{S}^\ell)_{ju}. \quad (14)$$

603 In particular, we denote by  $\mathbf{P}_{(vu)}^{(\ell)} \in \mathbb{R}^n$  the vector with entries  $(\mathbf{P}_{(vu)}^{(\ell)})_i = \mathbf{P}_{i(vu)}^{(\ell)}$ , for  $1 \leq i \leq n$ .

604 **Theorem C.2.** Given MPNNs as in (2), let  $\sigma$  and  $\psi^{(t)}$  be  $\mathcal{C}^2$  functions and assume  $|\sigma'|, |\sigma''| \leq c_\sigma$ ,  
 605  $\|\Omega^{(t)}\| \leq \omega$ ,  $\|\mathbf{W}^{(t)}\| \leq w$ ,  $\|\nabla_1 \psi^{(t)}\| \leq c_1$ ,  $\|\nabla_2 \psi^{(t)}\| \leq c_2$ ,  $\|\nabla^2 \psi^{(t)}\| \leq c^{(2)}$ . Let  $\mathbf{S} \in \mathbb{R}^{n \times n}$  be

$$\mathbf{S} := \frac{\omega}{w} \mathbf{I} + c_1 \text{diag}(\mathbf{A}\mathbf{1}) + c_2 \mathbf{A}.$$

606 Given nodes  $i, v, u \in \mathbf{V}$ , if  $\mathbf{P}_{(vu)}^{(\ell)} \in \mathbb{R}^n$  is as in (14) and  $m$  is the number of layers, then we derive

$$\begin{aligned} \|\nabla_{uv}^2 \mathbf{h}_i^{(m)}\| &\leq \sum_{k=0}^{m-1} \sum_{j \in \mathbf{V}} (c_\sigma w)^{2m-k-1} w (\mathbf{S}^{m-k})_{jv} (\mathbf{S}^k)_{ij} (\mathbf{S}^{m-k})_{ju} \\ &\quad + c^{(2)} \sum_{\ell=0}^{m-1} (c_\sigma w)^{m+\ell} (\mathbf{S}^{m-1-\ell} \mathbf{P}_{(vu)}^{(\ell)})_i. \end{aligned} \quad (15)$$

607 *Proof.* First, we note that  $\nabla_{uv}^2 \mathbf{h}_i^{(m)}$  is a matrix of dimension  $\mathbb{R}^{d \times (d \times d)}$ . We then use the following  
608 ordering for indexing the columns – which is consistent with a typical way of labelling columns of  
609 the Kronecker product of matrices, as detailed below (note that indices here start from 1):

$$\frac{\partial^2 h_i^{(m),\alpha}}{\partial x_v^\beta \partial x_u^\gamma} := \left( \nabla_{uv}^2 \mathbf{h}_i^{(m)} \right)_{\alpha, d(\beta-1)+\gamma}. \quad (16)$$

610 As above, we proceed by induction and start from the case  $m = 1$ :

$$\begin{aligned} \left( \nabla_{uv}^2 \mathbf{h}_i^{(1)} \right)_{\alpha, d(\beta-1)+\gamma} &= \sigma''(\tilde{h}_i^{(0),\alpha}) \left( \Omega_{\alpha\gamma}^{(1)} \delta_{iv} + \sum_r W_{\alpha r}^{(1)} \sum_j \mathbf{A}_{ij} (\delta_{iv} \partial_{1,\gamma} \psi^{(1),r} + \delta_{jv} \partial_{2,\gamma} \psi^{(1),r}) \right) \\ &\times \left( \Omega_{\alpha\beta}^{(1)} \delta_{iu} + \sum_r W_{\alpha r}^{(1)} \sum_j \mathbf{A}_{ij} (\delta_{iu} \partial_{1,\beta} \psi^{(1),r} + \delta_{ju} \partial_{2,\beta} \psi^{(1),r}) \right) \\ &+ \sigma'(\tilde{h}_i^{(0),\alpha}) \sum_r W_{\alpha r}^{(1)} \left( \sum_j \mathbf{A}_{ij} \delta_{iu} (\partial_{1,\gamma} \partial_{1,\beta} \psi^{(1),r} \delta_{iv} + \partial_{2,\gamma} \partial_{1,\beta} \psi^{(1),r} \delta_{jv}) \right) \\ &+ \sigma'(\tilde{h}_i^{(0),\alpha}) \sum_r W_{\alpha r}^{(1)} \left( \mathbf{A}_{iu} (\partial_{1,\gamma} \partial_{2,\beta} \psi^{(1),r} \delta_{iv} + \partial_{2,\gamma} \partial_{2,\beta} \psi^{(1),r} \delta_{uv}) \right) \\ &:= (Q_1)_{\alpha,\beta,\gamma} + (Q_2)_{\alpha,\beta,\gamma} + (Q_3)_{\alpha,\beta,\gamma}, \end{aligned}$$

611 where  $Q_1$  is the term containing second derivatives of  $\psi$  while  $Q_2, Q_3$  are the remaining expressions  
612 including second-order derivatives of the message functions  $\psi$ . Using the same strategy as for the  
613 first-order estimates, we can rewrite the first term  $Q_1$  as

$$\begin{aligned} (Q_1)_{\alpha,\beta,\gamma} &= \left( \text{diag}(\sigma''(\tilde{\mathbf{h}}_v^{(0)})) \left( \Omega^{(1)} \delta_{iv} + \mathbf{W}^{(1)} \left( \sum_j \mathbf{A}_{ij} \delta_{iv} \nabla_1 \psi^{(1)} + \mathbf{A}_{iv} \nabla_2 \psi^{(1)} \right) \right) \right)_{\alpha\gamma} \\ &\times \left( \Omega^{(1)} \delta_{iu} + \mathbf{W}^{(1)} \left( \sum_j \mathbf{A}_{ij} \delta_{iu} \nabla_1 \psi^{(1)} + \mathbf{A}_{iu} \nabla_2 \psi^{(1)} \right) \right)_{\alpha\beta} \end{aligned}$$

614 We now observe that given two matrices  $\mathbf{B}, \mathbf{C} \in \mathbb{R}^{d \times d}$  and  $1 \leq \alpha, \alpha', \beta, \gamma \leq d$ , the entries of the  
615 Kronecker product  $\mathbf{B} \otimes \mathbf{C}$  can be indexed as

$$(\mathbf{B} \otimes \mathbf{C})_{d(\alpha-1)+\alpha', d(\beta-1)+\gamma} = B_{\alpha\beta} C_{\alpha'\gamma}.$$

616 We now introduce the  $d \times (d \times d)$  sub-matrix of  $\mathbf{B} \otimes \mathbf{C}$  defined by

$$(\mathbf{B} \otimes \mathbf{C})'_{\alpha, d(\beta-1)+\gamma} = B_{\alpha\beta} C_{\alpha\gamma}. \quad (17)$$

617 Therefore, we can rewrite  $(Q_1)_{\alpha,\beta,\gamma}$  as the entry  $(\alpha, d(\beta-1) + \gamma)$  of the  $d \times (d \times d)$  sub-matrix

$$(\mathbf{Q}_1)_{\alpha, d(\beta-1)+\gamma} = (\mathbf{B} \otimes \mathbf{C})'_{\alpha, d(\beta-1)+\gamma}, \quad (18)$$

618 where

$$\begin{aligned} \mathbf{B} &:= \text{diag}(\sigma''(\tilde{\mathbf{h}}_v^{(0)})) \left( \Omega^{(1)} \delta_{iv} + \mathbf{W}^{(1)} \left( \sum_j \mathbf{A}_{ij} \delta_{iv} \nabla_1 \psi^{(1)} + \mathbf{A}_{iv} \nabla_2 \psi^{(1)} \right) \right), \\ \mathbf{C} &:= \Omega^{(1)} \delta_{iu} + \mathbf{W}^{(1)} \left( \sum_j \mathbf{A}_{ij} \delta_{iu} \nabla_1 \psi^{(1)} + \mathbf{A}_{iu} \nabla_2 \psi^{(1)} \right). \end{aligned}$$

619 Next, we proceed to write  $(Q_2)_{\alpha,\beta,\gamma}$  in matricial form. Before we do that, we observe that the  
620 Hessian of the message functions  $(\mathbf{x}_i, \mathbf{x}_j) \mapsto \psi^{(t)}(\mathbf{x}_i, \mathbf{x}_j)$  takes the form

$$\nabla^2 \psi^{(t)} = \begin{pmatrix} \nabla_{11}^2 \psi^{(t)} & \nabla_{12}^2 \psi^{(t)} \\ \nabla_{21}^2 \psi^{(t)} & \nabla_{22}^2 \psi^{(t)} \end{pmatrix},$$

621 where  $\nabla_{ab}^2 \psi^{(t)} \in \mathbb{R}^{d \times (d \times d)}$  and is indexed as follows

$$(\nabla_{ab}^2 \psi^{(t)})_{r, d(\beta-1)+\gamma} = \partial_{a,\beta} \partial_{b,\gamma} \psi^{(t),r},$$

622 where  $a, b \in \{1, 2\}$ . Using these notations, we note that

$$\sum_r W_{\alpha r}^{(1)} \partial_{1,\gamma} \partial_{1,\beta} \psi^{(1),r} = \left( \mathbf{W}^{(1)} \nabla_{11}^2 \psi^{(1)} \right)_{\alpha, d(\beta-1)+\gamma}.$$

623 Therefore, we derive

$$(Q_2)_{\alpha,\beta,\gamma} = (\mathbf{Q}_2)_{\alpha,d(\beta-1)+\gamma} = \sum_j \mathbf{A}_{ij} \delta_{iu} \delta_{iv} (\text{diag}(\sigma'(\tilde{\mathbf{h}}_i^{(0)})) \mathbf{W}^{(1)} \nabla_{11}^2 \psi^{(1)})_{\alpha,d(\beta-1)+\gamma} \\ + \mathbf{A}_{iv} \delta_{iu} (\text{diag}(\sigma'(\tilde{\mathbf{h}}_i^{(0)})) \mathbf{W}^{(1)} \nabla_{12}^2 \psi^{(1)})_{\alpha,d(\beta-1)+\gamma}. \quad (19)$$

624 A similar argument works for  $Q_3$ :

$$(Q_3)_{\alpha,\beta,\gamma} = (\mathbf{Q}_3)_{\alpha,d(\beta-1)+\gamma} = \mathbf{A}_{iu} \delta_{iv} (\text{diag}(\sigma'(\tilde{\mathbf{h}}_i^{(0)})) \mathbf{W}^{(1)} \nabla_{21}^2 \psi^{(1)})_{\alpha,d(\beta-1)+\gamma} \\ + \mathbf{A}_{iu} \delta_{uv} (\text{diag}(\sigma'(\tilde{\mathbf{h}}_i^{(0)})) \mathbf{W}^{(1)} \nabla_{22}^2 \psi^{(1)})_{\alpha,d(\beta-1)+\gamma}. \quad (20)$$

625 Therefore, we can combine (18), (19), and (20) to write

$$\|\nabla_{uv}^2 \mathbf{h}_i^{(1)}\| \leq \|\mathbf{Q}_1\| + \|\mathbf{Q}_2\| + \|\mathbf{Q}_3\| \\ \leq c_\sigma (\omega \delta_{iv} + w(c_1 \text{diag}(\mathbf{A}_1)_i \delta_{iv} + c_2 \mathbf{A}_{iv})) (\omega \delta_{iu} + w(c_1 \text{diag}(\mathbf{A}_1)_i \delta_{iu} + c_2 \mathbf{A}_{iu}) \\ + c_\sigma w c^{(2)} (\text{diag}(\mathbf{A}_1)_i \delta_{iv} \delta_{iu} + \mathbf{A}_{iv} \delta_{iu}) \\ + c_\sigma w c^{(2)} (\mathbf{A}_{iu} \delta_{iv} + \mathbf{A}_{iu} \delta_{uv})).$$

626 Finally, we can rely on (14) to re-arrange the equation above as

$$\|\nabla_{uv}^2 \mathbf{h}_i^{(1)}\| \leq (c_\sigma w) (w(\mathbf{S})_{iv} (\mathbf{S})_{iu}) + w c^{(2)} c_\sigma (\delta_{iv} \mathbf{A}_{iu} + \delta_{iu} \mathbf{A}_{iv} + \sum_j \delta_{jv} (\text{diag}(\mathbf{A}_1) + \mathbf{A})_{ij} \delta_{ju}) \\ = (c_\sigma w) (w(\mathbf{S})_{iv} (\mathbf{S})_{iu}) + c^{(2)} c_\sigma w \mathbf{P}_{i(vu)}^{(0)},$$

627 which proves the bound for the second-order derivatives in the case  $m = 1$ .

628 We now assume that the claim holds for all layers  $t \leq m - 1$ , and compute the second order derivative  
629 after  $m$  layers:

$$\left( \nabla_{uv}^2 \mathbf{h}_i^{(m)} \right)_{\alpha,d(\beta-1)+\gamma} = \sigma''(\tilde{h}_i^{(m-1),\alpha}) \\ \times \left( \sum_r \Omega_{\alpha r}^{(m)} (\nabla_u \mathbf{h}_i^{(m-1)})_{r\beta} \right. \\ \left. + W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{ij} \left( \sum_p \partial_{1,p} \psi^{(m),r} (\nabla_u \mathbf{h}_i^{(m-1)})_{p\beta} + \partial_{2,p} \psi^{(m),r} (\nabla_u \mathbf{h}_j^{(m-1)})_{p\beta} \right) \right) \\ \times \left( \sum_r \Omega_{\alpha r}^{(m)} (\nabla_v \mathbf{h}_i^{(m-1)})_{r\gamma} \right. \\ \left. + W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{ij} \left( \sum_p \partial_{1,p} \psi^{(m),r} (\nabla_v \mathbf{h}_i^{(m-1)})_{p\gamma} + \partial_{2,p} \psi^{(m),r} (\nabla_v \mathbf{h}_j^{(m-1)})_{p\gamma} \right) \right) \\ + \sigma'(\tilde{h}_i^{(m-1),\alpha}) \sum_r W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{ij} \sum_{p,q} \partial_{1,p} \partial_{1,q} \psi^{(m),r} (\nabla_u \mathbf{h}_i^{(m-1)})_{p\beta} (\nabla_v \mathbf{h}_i^{(m-1)})_{q\gamma} \\ + \sigma'(\tilde{h}_i^{(m-1),\alpha}) \sum_r W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{ij} \sum_{p,q} \partial_{1,p} \partial_{2,q} \psi^{(m),r} (\nabla_u \mathbf{h}_i^{(m-1)})_{p\beta} (\nabla_v \mathbf{h}_j^{(m-1)})_{q\gamma} \\ + \sigma'(\tilde{h}_i^{(m-1),\alpha}) \sum_r W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{ij} \sum_{p,q} \partial_{1,p} \partial_{2,q} \psi^{(m),r} (\nabla_u \mathbf{h}_j^{(m-1)})_{q\beta} (\nabla_v \mathbf{h}_i^{(m-1)})_{p\gamma} \\ + \sigma'(\tilde{h}_i^{(m-1),\alpha}) \sum_r W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{ij} \sum_{p,q} \partial_{2,p} \partial_{2,q} \psi^{(m),r} (\nabla_u \mathbf{h}_j^{(m-1)})_{p\beta} (\nabla_v \mathbf{h}_j^{(m-1)})_{q\gamma} \\ + \sigma'(\tilde{h}_i^{(m-1),\alpha}) \sum_r \Omega_{\alpha r}^{(m)} (\nabla_{uv}^2 \mathbf{h}_i^{(m-1)})_{r,d(\beta-1)+\gamma} \\ + \sigma'(\tilde{h}_i^{(m-1),\alpha}) \sum_r W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{ij} \left( \sum_p \partial_{1,p} \psi^{(m),r} (\nabla_{uv}^2 \mathbf{h}_i^{(m-1)})_{p,d(\beta-1)+\gamma} \right. \\ \left. + \sigma'(\tilde{h}_i^{(m-1),\alpha}) \sum_r W_{\alpha r}^{(m)} \sum_j \mathbf{A}_{ij} \left( \sum_p \partial_{2,p} \psi^{(m),r} (\nabla_{uv}^2 \mathbf{h}_j^{(m-1)})_{p,d(\beta-1)+\gamma} \right) \right) \\ := R_{\alpha,\beta,\gamma} + \sum_{a,b \in \{1,2\}} (Q_{ab})_{\alpha,\beta,\gamma} + Z_{\alpha,\beta,\gamma},$$

630 where  $R$  is the term containing second derivatives of the non-linear map  $\sigma$ ,  $Q_{ab}$  is indexed according  
 631 to the second derivatives of the message-functions  $\psi$ , and finally  $Z$  is the term containing second-  
 632 order derivatives of the features. For the term  $R_{\alpha,\beta,\gamma}$  we can argue as in the  $m = 1$  case and use the  
 633 sub-matrix notation in (17) to rewrite it as the entry  $(\alpha, d(\beta - 1) + \gamma)$  of the  $d \times (d \times d)$  sub-matrix

$$\mathbf{R}_{\alpha,d(\beta-1)+\gamma} = (\mathbf{B} \otimes \mathbf{C})'_{\alpha,d(\beta-1)+\gamma}, \quad (21)$$

634 where

$$\mathbf{B} := \text{diag}(\sigma''(\tilde{\mathbf{h}}_i^{(m-1)}))(\boldsymbol{\Omega}^{(m)} \nabla_u \mathbf{h}_i^{(m-1)} + \mathbf{W}^{(m)} (\sum_j \mathbf{A}_{ij} \nabla_1 \psi^{(m)} \nabla_u \mathbf{h}_i^{(m-1)} + \nabla_2 \psi^{(m)} \nabla_u \mathbf{h}_j^{(m-1)})),$$

$$\mathbf{C} := \boldsymbol{\Omega}^{(m)} \nabla_v \mathbf{h}_i^{(m-1)} + \mathbf{W}^{(m)} (\sum_j \mathbf{A}_{ij} \nabla_1 \psi^{(m)} \nabla_v \mathbf{h}_i^{(m-1)} + \nabla_2 \psi^{(m)} \nabla_v \mathbf{h}_j^{(m-1)})$$

635 Next we consider the terms  $(Q_{ab})_{\alpha,\beta,\gamma}$ . Without loss of generality, we focus on  $(Q_{11})_{\alpha,\beta,\gamma}$  and use  
 636 again the same argument in the  $m = 1$  case, to rewrite it as  $(Q_{11})_{\alpha,\beta,\gamma} = (\mathbf{Q}_{11})_{\alpha,d(\beta-1)+\gamma}$  where

$$\mathbf{Q}_{11} = \text{diag}(\sigma'(\tilde{\mathbf{h}}_i^{(m-1)})) \sum_j \mathbf{A}_{ij} (\mathbf{W}^{(m)} \nabla_{11}^2 \psi^{(m)} \nabla_u \mathbf{h}_i^{(m-1)} \otimes \nabla_v \mathbf{h}_i^{(m-1)}), \quad (22)$$

637 where again we are indexing the matrix  $\nabla_{11}^2 \psi^{(m)}$  by

$$(\nabla_{11}^2 \psi^{(m)})_{r,p(d-1)+q} = \partial_{1,p} \partial_{1,q} \psi^{(m),r}.$$

638 The other  $Q$ -terms can be estimated similarly. Finally, we rewrite  $Z_{\alpha,\beta,\gamma} = (\mathbf{Z})_{\alpha,d(\beta-1)+\gamma}$ , where

$$\mathbf{Z} = \text{diag}(\sigma'(\tilde{\mathbf{h}}_i^{(m-1)})) (\boldsymbol{\Omega}^{(m)} \nabla_{uv}^2 \mathbf{h}_i^{(m-1)} + \mathbf{W}^{(m)} \sum_j \mathbf{A}_{ij} (\nabla_1 \psi^{(m)} \nabla_{uv}^2 \mathbf{h}_i^{(m-1)} + \nabla_2 \psi^{(m)} \nabla_{uv}^2 \mathbf{h}_j^{(m-1)})) \quad (23)$$

639 Therefore, we have rewritten the second-derivatives of the features in matricial form as

$$\nabla_{uv}^2 \mathbf{h}_i^{(m)} = \mathbf{R} + \sum_{a,b \in \{1,2\}} \mathbf{Q}_{ab} + \mathbf{Z}.$$

640 To complete the proof, we now simply need to estimate the three terms and show they fit the recursion  
 641 claimed for  $m$ . For the case of  $\mathbf{R}$  in (21), we find

$$\|\mathbf{R}\| \leq c_\sigma (\omega \|\nabla_u \mathbf{h}_i^{(m-1)}\| + w(c_1 \text{diag}(\mathbf{A}\mathbf{1})_i \|\nabla_u \mathbf{h}_i^{(m-1)}\| + c_2 \sum_j \mathbf{A}_{ij} \|\nabla_u \mathbf{h}_j^{(m-1)}\|))$$

$$\times (\omega \|\nabla_v \mathbf{h}_i^{(m-1)}\| + w(c_1 \text{diag}(\mathbf{A}\mathbf{1})_i \|\nabla_v \mathbf{h}_i^{(m-1)}\| + c_2 \sum_j \mathbf{A}_{ij} \|\nabla_v \mathbf{h}_j^{(m-1)}\|).$$

642 If we write  $D\mathbf{h}^{(m-1)} \in \mathbb{R}^{n \times n}$  as the matrix with entries  $(D\mathbf{h}^{(m-1)})_{ij} = \|\nabla_j \mathbf{h}_i^{(m-1)}\|$ , then we  
 643 obtain

$$\|\mathbf{R}\| \leq c_\sigma w (\mathbf{w} \mathbf{S} D\mathbf{h}^{(m-1)})_{iv} (\mathbf{S} D\mathbf{h}^{(m-1)})_{iu}.$$

644 We can then plug the first-order estimates derived in Theorem C.1 and obtain

$$\|\mathbf{R}\| \leq c_\sigma w (\mathbf{w} \mathbf{S} (c_\sigma \mathbf{w})^{m-1} \mathbf{S}^{m-1})_{iv} (\mathbf{S} (c_\sigma \mathbf{w})^{m-1} \mathbf{S}^{m-1})_{iu} = (c_\sigma \mathbf{w})^{2m-1} (\mathbf{w} (\mathbf{S}^m)_{iv} (\mathbf{S}^m)_{iu}). \quad (24)$$

645 Next, we move onto the  $Q$ -terms, and use again the first-order estimates in Theorem C.1 – and the  
 646 fact that we can bound the norm of  $\nabla_{ab}^2 \psi^{(m)}$  by  $c^{(2)}$  – to derive

$$\|\sum_{a,b \in \{1,2\}} \mathbf{Q}_{ab}\| \leq c^{(2)} (c_\sigma \mathbf{w})^{2m-1} (\text{diag}(\mathbf{A}\mathbf{1})_i (\mathbf{S}^{m-1})_{iv} (\mathbf{S}^{m-1})_{iu} + \sum_j \mathbf{A}_{ij} (\mathbf{S}^{m-1})_{ju} (\mathbf{S}^{m-1})_{jv})$$

$$+ c^{(2)} (c_\sigma \mathbf{w})^{2m-1} ((\mathbf{S}^{m-1})_{iv} (\mathbf{A} \mathbf{S}^{m-1})_{iu} + (\mathbf{S}^{m-1})_{iu} (\mathbf{A} \mathbf{S}^{m-1})_{iv})$$

$$= c^{(2)} (c_\sigma \mathbf{w})^{2m-1} \mathbf{P}_{i(vu)}^{(m-1)}. \quad (25)$$

647 Finally, if we let  $D^2 \mathbf{h}_{vu} \in \mathbb{R}^n$  be the vector with entries  $(D^2 \mathbf{h}_{vu})_i = \|\nabla_{uv}^2 \mathbf{h}_i^{(m-1)}\|$ , then



$$\begin{aligned} \|\mathbf{Z}\| &\leq c_\sigma \left( \omega \|\nabla_{uv}^2 \mathbf{h}_i^{(m-1)}\| + w \left( c_1 \text{diag}(\mathbf{A}\mathbf{1})_i \|\nabla_{uv}^2 \mathbf{h}_i^{(m-1)}\| + c_2 \sum_j \mathbf{A}_{ij} \|\nabla_{uv}^2 \mathbf{h}_j^{(m-1)}\| \right) \right) \quad (26) \\ &= c_\sigma w (\mathbf{S} D^2 \mathbf{h}_{vu})_i. \end{aligned}$$

648 Therefore, we can use the induction to derive

$$\begin{aligned} \|\mathbf{Z}\| &\leq c_\sigma w \sum_s \mathbf{S}_{is} \sum_{k=0}^{m-2} \sum_{j \in \mathcal{V}} (c_\sigma w)^{2m-2-k-1} w (\mathbf{S}^{m-1-k})_{jv} (\mathbf{S}^k)_{sj} (\mathbf{S}^{m-1-k})_{ju} \\ &\quad + c_\sigma w \sum_s \mathbf{S}_{is} (c^{(2)}) \sum_{\ell=0}^{m-2} (c_\sigma w)^{m-1+\ell} (\mathbf{S}^{m-2-\ell} \mathbf{P}_{(vu)}^{(\ell)})_s \\ &= \sum_{k=0}^{m-2} \sum_{j \in \mathcal{V}} (c_\sigma w)^{2m-k-2} w (\mathbf{S}^{m-1-k})_{jv} (\mathbf{S}^{k+1})_{ij} (\mathbf{S}^{m-1-k})_{ju} \\ &\quad + c^{(2)} \left( \sum_{\ell=0}^{m-2} (c_\sigma w)^{m+\ell} (\mathbf{S}^{m-1-\ell} \mathbf{P}_{(vu)}^{(\ell)})_i \right) \\ &= \sum_{k=1}^{m-1} \sum_{j \in \mathcal{V}} (c_\sigma w)^{2m-k-1} w (\mathbf{S}^{m-k})_{jv} (\mathbf{S}^k)_{ij} (\mathbf{S}^{m-k})_{ju} + c^{(2)} \left( \sum_{\ell=0}^{m-2} (c_\sigma w)^{m+\ell} (\mathbf{S}^{m-1-\ell} \mathbf{P}_{(vu)}^{(\ell)})_i \right) \end{aligned}$$

649 By (24), we derive that the  $\mathbf{R}$ -term corresponds to the  $k = 0$  entry of the first sum, while (25)  
650 corresponds to the case  $\ell = m - 1$  of the second sum, which completes the induction and hence our  
651 proof.  $\square$

### 652 C.1 Proof of Theorem 3.2

653 We can now use the previous characterization to derive estimates on the Hessian of the graph-level  
654 function computed by MPNNs. We restate Theorem 3.2 here for convenience.

655 **Theorem 3.2.** Consider an MPNN of depth  $m$  as in (2), where  $\sigma$  and  $\psi^{(t)}$  are  $\mathcal{C}^2$  functions and we  
656 denote the bounds on their derivatives and on the norm of the weights as above. Let  $\mathbf{S}$  and  $\mathbf{Q}_k$  be  
657 defined as in (5) and (6), respectively. If the readout is MAX, MEAN or SUM and  $\theta$  in (3) has unit  
658 norm, then the mixing  $\text{mix}_{y_G^{(m)}}(v, u)$  induced by the MPNN over the features of nodes  $v, u$  satisfies

$$\text{mix}_{y_G^{(m)}}(v, u) \leq \sum_{k=0}^{m-1} (c_\sigma w)^{2m-k-1} \left( w (\mathbf{S}^{m-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) \mathbf{S}^{m-k} + c^{(2)} \mathbf{Q}_k \right)_{vu}. \quad (7)$$

659 *Proof.* First, we recall that according to Definition 3.1, we are interested in bounding the quantity

$$\text{mix}_{y_G^{(m)}}(v, u) = \max_{\mathbf{X}} \max_{1 \leq \beta, \gamma \leq d} \left| \frac{\partial^2 y_G^{(m)}(\mathbf{X})}{\partial x_u^\beta \partial x_v^\gamma} \right|.$$

660 Let us first consider the choice READ = SUM, so that by (3) we get

$$\text{mix}_{y_G^{(m)}}(v, u) \leq \left| \sum_{\alpha=1}^d \theta_\alpha \sum_{i \in \mathcal{V}} \frac{\partial^2 h_i^{(m), \alpha}}{\partial x_u^\beta \partial x_v^\gamma} \right|$$

661 As before, we index the columns of the Hessian of  $\mathbf{h}_i$  as  $\frac{\partial^2 h_i^{(m), \alpha}}{\partial x_u^\beta \partial x_v^\gamma} = (\nabla_{uv}^2 \mathbf{h}_i^{(m)})_{\alpha, d(\beta-1)+\gamma}$  and  
662 hence obtain

$$\text{mix}_{y_G^{(m)}}(v, u) \leq \sum_{i \in \mathcal{V}} \|(\nabla_{uv}^2 \mathbf{h}_i^{(m)})^\top \theta\| \leq \sum_{i \in \mathcal{V}} \|\nabla_{uv}^2 \mathbf{h}_i^{(m)}\|, \quad (27)$$

663 since  $\boldsymbol{\theta}$  has unit norm. Note that the very same bound in (27) also holds if we replaced the SUM  
 664 readout with either the MAX or the MEAN readout. We can then rely on Theorem C.2 and find

$$\begin{aligned}
 \text{mix}_{y_G}^{(m)}(v, u) &\leq \sum_{i \in V} \sum_{k=0}^{m-1} \sum_{j \in V} (c_\sigma \mathbf{w})^{2m-k-1} \mathbf{w}(\mathbf{S}^{m-k})_{jv} (\mathbf{S}^k)_{ij} (\mathbf{S}^{m-k})_{ju} \\
 &\quad + c^{(2)} \sum_{i \in V} \sum_{\ell=0}^{m-1} (c_\sigma \mathbf{w})^{m+\ell} (\mathbf{S}^{m-1-\ell} \mathbf{P}_{(vu)}^{(\ell)})_i \\
 &= \sum_{k=0}^{m-1} (c_\sigma \mathbf{w})^{2m-k-1} \left( \mathbf{w}(\mathbf{S}^{m-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) \mathbf{S}^{m-k} \right)_{vu} \\
 &\quad + c^{(2)} \sum_{\ell=0}^{m-1} (c_\sigma \mathbf{w})^{m+\ell} (\mathbf{1}^\top \mathbf{S}^{m-1-\ell}) \mathbf{P}_{(vu)}^{(\ell)} \\
 &= \sum_{k=0}^{m-1} (c_\sigma \mathbf{w})^{2m-k-1} \left( \mathbf{w}(\mathbf{S}^{m-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) \mathbf{S}^{m-k} \right)_{vu} \\
 &\quad + c^{(2)} \sum_{k=0}^{m-1} (c_\sigma \mathbf{w})^{2m-k-1} (\mathbf{1}^\top \mathbf{S}^k) \mathbf{P}_{(vu)}^{(m-k-1)}.
 \end{aligned}$$

665 For the second term we can simply use the formula in (14) to rewrite it in matricial form as  
 666 claimed – recall that  $\mathbf{Q}_k$  is defined in (6).  $\square$

## 667 D Proofs and additional details of Section 4

668 Throughout this Section, for simplicity, we assume  $c_\sigma = \max\{|c'_\sigma|, |c''_\sigma|\}$  to be smaller or equal than  
 669 one – this is satisfied by the vast majority of commonly used non-linear activations, and extending  
 670 the results below to arbitrary  $c_\sigma$  is straightforward.

### 671 D.1 Proof of Theorem 4.3

672 We begin by proving lower bounds on the operator norm of the weights, when the depth is the minimal  
 673 one required to induce any non-zero mixing among nodes  $v, u$ . For convenience, we restate Theorem  
 674 4.3 here as well.

675 **Theorem 4.3.** Let  $\mathbf{A} = \mathbf{A}_{\text{sym}}$ ,  $r := d_G(v, u)$ ,  $m = \lceil r/2 \rceil$ , and  $q$  be the number of paths of length  
 676  $r$  between  $v$  and  $u$ . For an MPNN satisfying Theorem 3.2 with capacity  $(m = \lceil r/2 \rceil, \mathbf{w})$ , we find  
 677  $\widetilde{\text{OSQ}}_{v,u}(m, \mathbf{w}) \cdot (c_2 \mathbf{w})^r (\mathbf{A}^r)_{vu} \geq 1$ . In particular, if the MPNN generates mixing  $\text{mix}_{y_G}(v, u)$ , then

$$\mathbf{w} \geq \frac{d_{\min}}{c_2} \left( \frac{\text{mix}_{y_G}(v, u)}{q} \right)^{\frac{1}{r}}.$$

678 *Proof.* Without loss of generality, we assume that  $r$  is even, so that by our assumptions, we can  
 679 simply take  $m = r/2$ . According to Theorem 3.2, we know that the maximal mixing induced by an  
 680 MPNN of depth  $m$  over the features associated with nodes  $v, u$  is bounded from above as

$$\text{mix}_{y_G}^{(m)}(v, u) \leq \sum_{k=0}^{m-1} \mathbf{w}^{2m-k-1} \left( \mathbf{w} \mathbf{S}^{m-k} \text{diag}(\mathbf{S}^k \mathbf{1}) \mathbf{S}^{m-k} + c^{(2)} \mathbf{Q}_k \right)_{vu}.$$

681 where we have replaced  $c_\sigma$  with one, as per our assumption. Since  $m = r/2$ , where  $r$  is the  
 682 distance among nodes  $v, u$ , then the only non-zero contribution for the first term is obtained for  
 683  $k = 0$  – otherwise we would find a path of length  $2(m - k)$  connecting  $v$  and  $u$  hence violating the  
 684 assumptions – and is equal to  $\mathbf{w}^{2m} (\mathbf{S}^{2m})_{vu}$ , and note that  $2m = r$ . Concerning the terms  $\mathbf{Q}_k$  instead,  
 685 the longest-walk contribution for nodes  $v, u$  is  $2m - 1$  (when  $k = 0$ ), meaning that  $\mathbf{Q}_k = 0$  for all  
 686  $0 \leq k \leq m - 1$  if  $m = 2r$ . Accordingly, we can reduce the bound above to:

$$\text{mix}_{y_G}^{(m)}(v, u) \leq \mathbf{w}^{2m} (\mathbf{S}^{2m})_{vu} = \mathbf{w}^{2m} \left( (c_2 \mathbf{A})^{2m} \right)_{vu},$$

687 where in the last equality we have again used that  $2m = r$ , so when expanding the power of  $\mathbf{S}$  only  
 688 the highest-order term in the  $\mathbf{A}$ -variable gives non-zero contributions. If we replace now  $2m = r$  and  
 689 use the characterization of over-squashing in Definition 4.2, then

$$\widetilde{\text{OSQ}}_{v,u}(m = \frac{r}{2}, \mathbf{w}) \geq (c_2 \mathbf{w})^{-r} \frac{1}{(\mathbf{A}^r)_{vu}}.$$

690 Therefore, if the MPNN generates mixing  $\text{mix}_{y_G^{(m)}}(v, u)$  among the features of  $v$  and  $u$ , then (8) is  
 691 satisfied, meaning that the operator norm of the weights must be larger than

$$\mathbf{w} \geq \frac{1}{c_2} \left( \frac{\text{mix}_{y_G^{(m)}}(v, u)}{(\mathbf{A}^r)_{vu}} \right)^{\frac{1}{r}}.$$

692 The term  $\mathbf{A}^r$  in general can be estimated sharply depending on the knowledge we have of the  
 693 underlying graph. To get a universal – albeit potentially looser bound – it suffices to note that along  
 694 each path connecting  $v$  and  $u$ , the product of the entries of  $\mathbf{A}$  can be bounded from above by  $(d_{\min})^r$ ,  
 695 which completes the proof.  $\square$

696 We highlight that if  $\mathbf{A} = \mathbf{A}_{\text{rw}} = \mathbf{D}^{-1} \mathbf{A}$  – i.e. the aggregation over the neighbours consists of a  
 697 mean-operation as for the GraphSAGE architecture – then one can apply the very same proof above  
 698 and derive

699 **Corollary D.1.** *The same lower bound for  $w$  in Theorem 4.3, holds when  $\mathbf{A} = \mathbf{D}^{-1} \mathbf{A}$ .*

700 **Some examples.** We illustrate the bounds in Theorem 4.3 and for simplicity, we set  $c_2 = 1$ .  
 701 Consider a tree  $T_d$  of arity  $d$ , with  $v$  the root and  $u$  a leaf at distance  $r$  and depth  $m = r/2$ ; then  
 702  $\widetilde{\text{OSQ}}_{v,u}(m, \mathbf{w}) \geq w^{-r} (d+1)^{r-1}$  and the operator norm required to generate mixing  $y(v, u)$  is

$$\mathbf{w} \geq (d+1) \left( \frac{y(v, u)}{d+1} \right)^{\frac{1}{r}}.$$

703 We note that by taking  $d = 1$  we recover the case of the path-graph (1D grid). Since the operator  
 704 norm of the weights grows with the branching factor, we see that, in general, the capacity required by  
 705 MPNNs to solve long-range tasks could be higher on graphs than on sequences [2]. We also consider  
 706 the case of a 1-layer MPNN on a complete graph  $K_n$  with  $v \neq u$ . We find that  $\widetilde{\text{OSQ}}_{v,u}(m, \mathbf{w}) \geq$   
 707  $(n-1)/w$  and hence the operator norm required to generate mixing  $y(v, u)$  is  $w \geq (n-1)y(v, u)$ .  
 708 We note how the measure of over-squashing also captures the problem of redundancy of messages  
 709 [14]. In fact, even if  $v, u$  are at distance 1, the more nodes are there in the complete graph and hence  
 710 the more messages are exchanged, the more difficult for a shallow MPNN to induce enough mixing  
 711 among those *specific* nodes.

## 712 D.2 Spectral bounds

713 Next, we study the case of fixed, bounded operator norm of the weights, but variable depth, since  
 714 we are interested in showing that *over-squashing hinders the expressive power of MPNNs for tasks*  
 715 *requiring high-mixing of features associated with nodes at high commute time*. We first provide a  
 716 characterization of the maximal mixing (and hence of the over-squashing measure) in terms of the  
 717 graph-Laplacian and its pseudo-inverse.

718 **Convention.** In the proofs below we usually deal with matrices with nonnegative entries. Accordingly,  
 719 we introduce the following convention: we write that  $\mathbf{A} \leq \mathbf{B}$  if  $A_{ij} \leq B_{ij}$  for all entries  $1 \leq i, j \leq n$ .

720 **Theorem D.2.** *Let  $\gamma := \sqrt{\frac{d_{\max}}{d_{\min}}}$  and set  $\mathbf{A} = \mathbf{A}_{\text{sym}}$  or  $\mathbf{A} = \mathbf{A}_{\text{rw}}$ . Consider an MPNN as in Thm. 3.2  
 721 with depth  $m$ ,  $\max\{w, \omega/w + c_1 \gamma + c_2\} \leq 1$ . Define  $\mathbf{Z} := \mathbf{I} - c_2 \mathbf{\Delta}$ . Then the maximal mixing of  
 722 nodes  $v, u$  generated by such MPNN after  $m$  layers is*

$$\begin{aligned} \text{mix}_{y_G^{(m)}}(v, u) &\leq \gamma^k \left( m \frac{\sqrt{d_v d_u}}{2|E|} \left( 1 + 2c^{(2)}(1 + \gamma^s) \right) + \frac{1}{c_2} \left( \mathbf{Z}^2 (\mathbf{I} - \mathbf{Z}^{2m}) (\mathbf{I} + \mathbf{Z})^{-1} \mathbf{\Delta}^\dagger \right)_{vu} \right) \\ &\quad + 2 \frac{c^{(2)}}{c_2} \gamma^k \left( \left( (1 + \gamma^s) \mathbf{I} - \mathbf{\Delta} \right) (\mathbf{I} - \mathbf{Z}^{2m}) (\mathbf{I} + \mathbf{Z})^{-1} \mathbf{\Delta}^\dagger \right)_{vu}, \end{aligned}$$

723 where  $k = s = 1$  if  $\mathbf{A} = \mathbf{A}_{\text{sym}}$  or  $k = 4, s = 2$  if  $\mathbf{A} = \mathbf{A}_{\text{rw}}$ .

724 *Proof.* We first focus on the symmetrically normalized case  $\mathbf{A} = \mathbf{A}_{\text{sym}} = \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ , which  
725 we recall that we can rewrite as  $\mathbf{A}_{\text{sym}} = \mathbf{I} - \mathbf{\Delta}$ , where  $\mathbf{\Delta}$  is the (normalized) graph Laplacian (9).  
726 Since  $w \leq 1$  and the message-passing matrix is symmetric, by Theorem 3.2, we can bound the  
727 maximal mixing of an MPNN as in the statement by

$$\text{mix}_{y_G^{(m)}}(v, u) \leq \sum_{k=0}^{m-1} \left( \mathbf{S}^{m-k} \text{diag}(\mathbf{S}^k \mathbf{1}) \mathbf{S}^{m-k} + c^{(2)} \mathbf{Q}_k \right)_{vu}.$$

728 We focus on the first sum. Since  $\mathbf{A} = \mathbf{D}^{1/2}(\mathbf{D}^{-1}\mathbf{A})\mathbf{D}^{-1/2}$ , where  $\mathbf{D}^{-1}\mathbf{A}$  is a row-stochastic matrix,  
729 we see that

$$S_{ij} \leq \frac{\omega}{w} \delta_{ij} + c_1 \gamma \delta_{ij} + c_2 A_{ij},$$

730 meaning that we can write  $\mathbf{S} \leq \alpha \mathbf{I} + c_2 \mathbf{A}$ , where  $\alpha = \omega/w + c_1 \gamma$ , using the convention introduced  
731 above. Accordingly, we can estimate the row-sum of the powers of  $\mathbf{S}$  by using  $(\mathbf{A}^p \mathbf{1})_i \leq \gamma$  as

$$(\mathbf{S}^k \mathbf{1})_i \leq \sum_{p=0}^k \binom{k}{p} \alpha^{k-p} c_2^p (\mathbf{A}^p \mathbf{1})_i \leq \gamma \sum_{p=0}^k \binom{k}{p} \alpha^{k-p} c_2^p = \gamma (\alpha + c_2)^k \leq \gamma,$$

732 where the last inequality simply follows from the assumptions. Therefore, we find

$$\left( \sum_{k=0}^{m-1} \mathbf{S}^{m-k} \text{diag}(\mathbf{S}^k \mathbf{1}) \mathbf{S}^{m-k} \right)_{vu} \leq \gamma \sum_{k=0}^{m-1} (\mathbf{S}^{2(m-k)})_{vu} = \gamma \sum_{k=1}^m (\mathbf{S}^{2k})_{vu}. \quad (28)$$

733 By the assumptions on the regularity of the message-functions, we can estimate  $\mathbf{S}$  from above by  
734  $\mathbf{S} \leq \alpha \mathbf{I} + c_2 \mathbf{A} = (\alpha + c_2) \mathbf{I} - c_2 \mathbf{\Delta} \leq \mathbf{Z}$ , and derive

$$\left( \sum_{k=0}^{m-1} \mathbf{S}^{m-k} \text{diag}(\mathbf{S}^k \mathbf{1}) \mathbf{S}^{m-k} \right)_{vu} \leq \gamma \sum_{k=1}^m (\mathbf{Z}^{2k})_{vu}.$$

735 From the spectral decomposition of the graph-Laplacian in (9) and the properties that  $\lambda_0 = 0$  and  
736  $\phi_0(v) = \sqrt{d_v/2|E|}$ , we find

$$\begin{aligned} \sum_{k=1}^m (\mathbf{Z}^{2k})_{vu} &= \sum_{k=1}^m \sum_{\ell=0}^{n-1} (1 - c_2 \lambda_\ell)^{2k} \phi_\ell(v) \phi_\ell(u) \\ &= m \frac{\sqrt{d_v d_u}}{2|E|} + \sum_{\ell=1}^{n-1} \left( \frac{1 - (1 - c_2 \lambda_\ell)^{2(m+1)}}{1 - (1 - c_2 \lambda_\ell)^2} - 1 \right) \phi_\ell(v) \phi_\ell(u) \\ &= m \frac{\sqrt{d_v d_u}}{2|E|} + \sum_{\ell=1}^{n-1} \left( \frac{(1 - c_2 \lambda_\ell)^2 (1 - (1 - c_2 \lambda_\ell)^{2m})}{(2 - c_2 \lambda_\ell) c_2 \lambda_\ell} \right) \phi_\ell(v) \phi_\ell(u). \end{aligned}$$

737 Since  $c_2 \leq 1$  and  $G$  is not bipartite, we derive that  $(\mathbf{I} + \mathbf{Z}) = 2\mathbf{I} - c_2 \mathbf{\Delta}$  is invertible and hence that  
738 the following decomposition holds:

$$(\mathbf{I} + \mathbf{Z})^{-1} = \sum_{\ell \geq 0} \frac{1}{2 - c_2 \lambda_\ell} \phi_\ell \phi_\ell^\top.$$

739 Therefore, we can rely on the spectral-decomposition of the pseudo-inverse of the graph-Laplacian in  
740 (10) to get

$$\left( \sum_{k=0}^{m-1} \mathbf{S}^{m-k} \text{diag}(\mathbf{S}^k \mathbf{1}) \mathbf{S}^{m-k} \right)_{vu} \leq \gamma \left( m \frac{\sqrt{d_v d_u}}{2|E|} + \frac{1}{c_2} \left( \mathbf{Z}^2 (\mathbf{I} - \mathbf{Z}^{2m}) (\mathbf{I} + \mathbf{Z})^{-1} \mathbf{\Delta}^\dagger \right)_{vu} \right). \quad (29)$$

741 It now remains to bound the term  $c^{(2)} \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu}$ . First, we note that by the symmetry of  $\mathbf{A}$  and  
742 the estimate  $(\mathbf{S}^k \mathbf{1})_i \leq \gamma$ , that we derived above, we obtain

$$c^{(2)} \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu} \leq 2c^{(2)} \gamma \left( \sum_{k=0}^{m-1} (\mathbf{A} \mathbf{Z}^{2(m-k-1)})_{vu} + \gamma (\mathbf{Z}^{2(m-k-1)})_{vu} \right).$$

743 Then we can use the identity  $\mathbf{A} = \mathbf{I} - \Delta$ , to find

$$c^{(2)} \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu} \leq 2c^{(2)} \gamma \sum_{k=0}^{m-1} \left( ((1 + \gamma)\mathbf{I} - \Delta) \mathbf{Z}^{2(m-k-1)} \right)_{vu}. \quad (30)$$

744 By relying on the spectral decomposition as above, we finally get

$$\begin{aligned} c^{(2)} \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu} &\leq 2c^{(2)} \gamma \left( m \frac{\sqrt{d_v d_u}}{2|\mathbf{E}|} (1 + \gamma) \right) \\ &\quad + 2c^{(2)} \gamma \left( \sum_{\ell > 0} (1 + \gamma - \lambda_\ell) \frac{1 - (1 - c_2 \lambda_\ell)^{2m}}{(2 - c_2 \lambda_\ell) c_2 \lambda_\ell} \phi_\ell(v) \phi_\ell(u) \right). \end{aligned}$$

745 As done previously, we can rewrite the last terms via  $(\mathbf{I} + \mathbf{Z})^{-1}$  as

$$\begin{aligned} c^{(2)} \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu} &\leq 2c^{(2)} \gamma \left( m \frac{\sqrt{d_v d_u}}{2|\mathbf{E}|} (1 + \gamma) \right) \\ &\quad + 2 \frac{c^{(2)}}{c_2} \gamma \left( (1 + \gamma) \mathbf{I} - \Delta \right) (\mathbf{I} - \mathbf{Z}^{2m}) (\mathbf{I} + \mathbf{Z})^{-1} \Delta^\dagger \Big|_{vu}. \end{aligned} \quad (31)$$

746 We can then combine (29) and (31) and derive the bound we claimed, when  $\mathbf{A} = \mathbf{A}_{\text{sym}}$ . For the case  
747  $\mathbf{A} = \mathbf{A}_{\text{rw}} = \mathbf{D}^{-1} \mathbf{A}$ , it suffices to notice that  $\mathbf{S} \leq \alpha' \mathbf{I} + c_2 \mathbf{A}$ , where  $\alpha' = \omega/w + c_1$  and that

$$(\mathbf{1}^\top \mathbf{S}^k)_i \leq \sum_j \sum_{p=0}^k \binom{k}{p} (\alpha')^{k-p} c_2^p ((\mathbf{D}^{-1} \mathbf{A})^p)_{ji} \leq \frac{d_{\max}}{d_{\min}} (\alpha' + c_2)^k \leq \gamma^2,$$

748 where we have used that by assumption  $\alpha' + c_2 \leq 1$ . Similarly, we get  $(\mathbf{S}^{m-k})^\top \mathbf{S}^{(m-k)} \leq$   
749  $\gamma^2 \mathbf{Z}^{2(m-k)}$ . Finally, the  $\mathbf{Q}_k$ -term can be bounded by

$$c^{(2)} \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu} \leq 2c^{(2)} \left( \sum_{k=0}^{m-1} \gamma^4 \mathbf{A}_{\text{sym}} \mathbf{Z}^{2(m-k-1)} + \gamma^6 \mathbf{Z}^{2(m-k-1)} \right),$$

750 and we can follow the previous steps in the symmetric case to complete the proof.  $\square$

751 **Corollary D.3.** *Under the assumptions of Theorem D.2, if the message functions in (2) are linear –*  
752 *as for GCN, SAGE, or GIN – then the maximal mixing induced by such an MPNN of  $m$  layers is*

$$\text{mix}_{y_G^{(m)}}(v, u) \leq \left( \frac{d_{\max}}{d_{\min}} \right)^k \left( m \frac{\sqrt{d_v d_u}}{2|\mathbf{E}|} + \frac{1}{c_2} \left( \mathbf{Z}^2 (\mathbf{I} - \mathbf{Z}^{2m}) (\mathbf{I} + \mathbf{Z})^{-1} \Delta^\dagger \right)_{vu} \right)$$

753 *Proof.* This follows from Theorem D.2 simply by noticing that if the message-function  $\psi$  in (2) is  
754 linear, then the upper bound for the norm of the Hessian can be taken to be zero, i.e.  $c^{(2)} = 0$ .  $\square$

### 755 D.3 Proof of Theorem 4.4

756 We now expand the previous results to derive the minimal number of layers required to induce mixing  
757 in the case of bounded weights, showing that the depth may need to grow with the commute time of  
758 nodes. We recall that  $\gamma$  is  $\sqrt{d_{\max}/d_{\min}}$  while  $0 = \lambda_0 < \lambda_1 \leq \dots \leq \lambda_{n-1}$  are the eigenvalues of the  
759 symmetrically normalized graph Laplacian (9). We restate Theorem 4.4 below.

760 **Theorem 4.4.** *Consider an MPNN satisfying Theorem 3.2, with  $\max\{w, \omega/w + c_1 \gamma + c_2\} \leq 1$ ,*  
761 *and  $\mathbf{A} = \mathbf{A}_{\text{sym}}$ . If  $\widehat{\text{OSQ}}_{v,u}(m, w) \cdot (\text{mix}_{y_G}(v, u)) \leq 1$ , i.e. the MPNN generates mixing  $\text{mix}_{y_G}(v, u)$*   
762 *among the features associated with nodes  $v, u$ , then the number of layers  $m$  satisfies*

$$m \geq \frac{\tau(v, u)}{4c_2} + \frac{|\mathbf{E}|}{\sqrt{d_v d_u}} \left( \frac{\text{mix}_{y_G}(v, u)}{\gamma \mu} - \frac{1}{c_2} \left( \frac{\gamma + |1 - c_2 \lambda^*|^{r-1}}{\lambda_1} + 2 \frac{c^{(2)}}{\mu} \right) \right),$$

763 where  $r = d_G(v, u)$ ,  $\mu = 1 + 2c^{(2)}(1 + \gamma)$  and  $|1 - c_2 \lambda^*| = \max_{0 < \ell \leq n-1} |1 - c_2 \lambda_\ell| < 1$ .



764 *Proof.* From now on we let  $r$  be the shortest-walk distance between  $v$  and  $u$ . If  $m < r/2$ , then we  
765 incur the under-reaching issue and hence we get zero mixing among the features associated with  
766 nodes  $v, u$ . Accordingly, we can choose  $m \geq r/2$ . We need to provide an estimate on the maximal  
767 mixing induced by an MPNN as in the statement. We focus on the bound in Theorem 3.2, and recall  
768 that the first sum can be bounded as in (28) by

$$\left( \sum_{k=0}^{m-1} \mathbf{S}^{m-k} \text{diag}(\mathbf{S}^k \mathbf{1}) \mathbf{S}^{m-k} \right)_{vu} \leq \gamma \sum_{k=1}^m (\mathbf{S}^{2k})_{vu} \leq \gamma \sum_{k=1}^m (\mathbf{Z}^{2k})_{vu},$$

769 where  $\mathbf{Z} := \mathbf{I} - c_2 \mathbf{\Delta}$ . We can then bound the geometric sum by accounting for the odd powers too.  
770 Therefore, we get

$$\gamma \sum_{k=1}^m (\mathbf{Z}^{2k})_{vu} \leq \gamma \sum_{k=0}^{2m} (\mathbf{Z}^k)_{vu} = \gamma \sum_{k=0}^{2m} \sum_{\ell \geq 0} (1 - c_2 \lambda_\ell)^k \phi_\ell(v) \phi_\ell(u).$$

771 As for the proof of Theorem D.2, we separate the contribution of the zero-eigenvalue and that of the  
772 positive ones, so we find that

$$\begin{aligned} \sum_{k=0}^{2m} (\mathbf{Z}^k)_{vu} &\leq (2m+1) \frac{\sqrt{d_v d_u}}{2|\mathbf{E}|} + \sum_{\ell > 0} \left( \frac{1 - (1 - c_2 \lambda_\ell)^{2m+1}}{c_2 \lambda_\ell} \right) \phi_\ell(v) \phi_\ell(u) \\ &= (2m+1) \frac{\sqrt{d_v d_u}}{2|\mathbf{E}|} + \sum_{\ell > 0} \frac{1}{c_2 \lambda_\ell} \phi_\ell(v) \phi_\ell(u) - \sum_{\ell > 0} \frac{(1 - c_2 \lambda_\ell)^{2m+1}}{c_2 \lambda_\ell} \phi_\ell(v) \phi_\ell(u). \end{aligned} \quad (32)$$

773 Thanks to the characterization of commute-time provided in (11), we derive

$$\begin{aligned} \sum_{\ell > 0} \frac{1}{c_2 \lambda_\ell} \phi_\ell(v) \phi_\ell(u) &= -\frac{\tau(v, u)}{4c_2 |\mathbf{E}|} \sqrt{d_v d_u} + \frac{1}{2c_2} \sum_{\ell > 0} \frac{1}{\lambda_\ell} \left( \phi_\ell^2(v) \sqrt{\frac{d_u}{d_v}} + \phi_\ell^2(u) \sqrt{\frac{d_v}{d_u}} \right) \\ &\leq -\frac{\tau(v, u)}{4c_2 |\mathbf{E}|} \sqrt{d_v d_u} + \frac{1}{2c_2 \lambda_1} \left( \sqrt{\frac{d_v}{d_u}} + \sqrt{\frac{d_u}{d_v}} - \frac{\sqrt{d_v d_u}}{|\mathbf{E}|} \right) \end{aligned} \quad (33)$$

774 where in the last inequality we have used that  $\sum_{\ell > 0} \phi_\ell^2(v) = 1 - \phi_0^2(v)$  since  $\{\phi_\ell\}$  constitute an  
775 orthonormal basis, with  $\phi_0(v) = \sqrt{d_v/2|\mathbf{E}|}$ , and that  $\lambda_\ell \geq \lambda_1$ , for all  $\ell > 0$ . Next, we estimate the  
776 second sum in (32), and we note that  $\lambda^*$  in the statement is either  $\lambda_1$  or  $\lambda_{n-1}$ :

$$\begin{aligned} -\sum_{\ell > 0} \frac{(1 - c_2 \lambda_\ell)^{2m+1}}{c_2 \lambda_\ell} \phi_\ell(v) \phi_\ell(u) &\leq \sum_{\ell > 0} \frac{|1 - c_2 \lambda^*|^{2m+1}}{c_2 \lambda_\ell} |\phi_\ell(v) \phi_\ell(u)| \\ &\leq \frac{|1 - c_2 \lambda^*|^{2m+1}}{2c_2 \lambda_1} \sum_{\ell > 0} (\phi_\ell^2(v) + \phi_\ell^2(u)) \\ &\leq \frac{|1 - c_2 \lambda^*|^r}{2c_2 \lambda_1} \left( 2 - \frac{d_v}{2|\mathbf{E}|} - \frac{d_u}{2|\mathbf{E}|} \right), \end{aligned} \quad (34)$$

777 where in the last inequality we have used that  $|1 - c_2 \lambda^*| < 1$  and that  $m \geq r/2$  (otherwise we would  
778 have zero-mixing due to under-reaching). Therefore, by combining (33) and (34), we derive that the  
779 first sum on the right hand side of (7) can be bounded from above by

$$\begin{aligned} \left( \sum_{k=0}^{m-1} \mathbf{S}^{m-k} \text{diag}(\mathbf{S}^k \mathbf{1}) \mathbf{S}^{m-k} \right)_{vu} &\leq \gamma \left( (2m+1) \frac{\sqrt{d_v d_u}}{2|\mathbf{E}|} - \frac{\tau(v, u)}{4c_2 |\mathbf{E}|} \sqrt{d_v d_u} \right) \\ &\quad + \frac{\gamma}{2c_2 \lambda_1} \left( \sqrt{\frac{d_v}{d_u}} + \sqrt{\frac{d_u}{d_v}} - \frac{\sqrt{d_v d_u}}{|\mathbf{E}|} \right) \\ &\quad + \gamma \frac{|1 - c_2 \lambda^*|^r}{2c_2 \lambda_1} \left( 2 - \frac{d_v}{2|\mathbf{E}|} - \frac{d_u}{2|\mathbf{E}|} \right). \end{aligned} \quad (35)$$

780 Next, we continue by estimating the second sum entering the right hand side of (7). We recall that by  
 781 (30), we have

$$\begin{aligned}
 c^{(2)} \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu} &\leq 2c^{(2)}\gamma \sum_{k=0}^{m-1} \left( ((1+\gamma)\mathbf{I} - \Delta) \mathbf{Z}^{2(m-k-1)} \right)_{vu} = 2c^{(2)}\gamma \sum_{k=0}^{m-1} \left( ((1+\gamma)\mathbf{I} - \Delta) \mathbf{Z}^{2k} \right)_{vu} \\
 &\leq 2c^{(2)}\gamma \sum_{k=0}^{2(m-1)} \left( ((1+\gamma)\mathbf{I} - \Delta) \mathbf{Z}^k \right)_{vu} \\
 &= 2c^{(2)}\gamma \sum_{k=0}^{2(m-1)} \sum_{\ell \geq 0} ((1+\gamma) - \lambda_\ell) (1 - c_2 \lambda_\ell)^k \phi_\ell(v) \phi_\ell(u).
 \end{aligned}$$

782 We then proceed as above, and separate the contributions associated with the kernel of the Laplacian,  
 783 to find

$$\begin{aligned}
 \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu} &\leq 2\gamma \left( (1+\gamma)(2m-1) \frac{\sqrt{d_v d_u}}{2|\mathbf{E}|} \right) \\
 &\quad + 2\gamma(1+\gamma) \sum_{\ell > 0} \frac{1}{c_2 \lambda_\ell} (1 - (1 - c_2 \lambda_\ell)^{2m-1}) \phi_\ell(v) \phi_\ell(u) \tag{36} \\
 &\quad - \frac{2\gamma}{c_2} \sum_{\ell > 0} (1 - (1 - c_2 \lambda_\ell)^{2m-1}) \phi_\ell(v) \phi_\ell(u). \tag{37}
 \end{aligned}$$

784 For the term in (36), we can apply the same estimate as for the case of (32). Similarly, we can bound  
 785 (37) by

$$-\frac{2\gamma}{c_2} \sum_{\ell > 0} (1 - (1 - c_2 \lambda_\ell)^{2m-1}) \phi_\ell(v) \phi_\ell(u) \leq \frac{2\gamma}{c_2} \sum_{\ell > 0} \frac{1}{2} (\phi_\ell^2(v) + \phi_\ell^2(u)) \leq \frac{\gamma}{c_2} \left( 2 - \frac{d_v}{2|\mathbf{E}|} - \frac{d_u}{2|\mathbf{E}|} \right).$$

786 Therefore, we can finally bound the  $\mathbf{Q}_k$ -terms in (7) by

$$\begin{aligned}
 c^{(2)} \sum_{k=0}^{m-1} (\mathbf{Q}_k)_{vu} &\leq 2c^{(2)}\gamma \left( (1+\gamma)(2m-1) \frac{\sqrt{d_v d_u}}{2|\mathbf{E}|} - (1+\gamma) \frac{\tau(v, u)}{4c_2 |\mathbf{E}|} \sqrt{d_v d_u} \right) \\
 &\quad + c^{(2)}\gamma \left( \frac{1+\gamma}{c_2 \lambda_1} \left( \sqrt{\frac{d_v}{d_u}} + \sqrt{\frac{d_u}{d_v}} - \frac{\sqrt{d_v d_u}}{|\mathbf{E}|} \right) \right) \\
 &\quad + c^{(2)}\gamma \frac{1+\gamma}{c_2 \lambda_1} |1 - c_2 \lambda^*|^{r-1} \left( 2 - \frac{d_v}{2|\mathbf{E}|} - \frac{d_u}{2|\mathbf{E}|} \right) \\
 &\quad + \frac{c^{(2)}\gamma}{c_2} \left( 2 - \frac{d_v}{2|\mathbf{E}|} - \frac{d_u}{2|\mathbf{E}|} \right). \tag{38}
 \end{aligned}$$

787 We can the combine (35) and (38), to find that the maximal mixing induced by an MPNN of  $m$  layers  
 788 as in the statement of Theorem 4.4, is

$$\begin{aligned}
 \text{mix}_{y_G^{(m)}}(v, u) &\leq \gamma \sqrt{d_v d_u} \left( \frac{m}{|\mathbf{E}|} \mu + \frac{1}{2|\mathbf{E}|} - \frac{\tau(v, u)}{4c_2 |\mathbf{E}|} \mu \right) \\
 &\quad + \gamma \frac{\mu}{2c_2 \lambda_1} \left( \sqrt{\frac{d_v}{d_u}} + \sqrt{\frac{d_u}{d_v}} - \frac{\sqrt{d_v d_u}}{|\mathbf{E}|} \right) \\
 &\quad + \gamma \frac{\mu}{2c_2 \lambda_1} |1 - c_2 \lambda^*|^{r-1} \left( 2 - \frac{d_v}{2|\mathbf{E}|} - \frac{d_u}{2|\mathbf{E}|} \right) \\
 &\quad + \frac{\gamma c^{(2)}}{c_2} \left( 2 - \frac{d_v}{2|\mathbf{E}|} - \frac{d_u}{2|\mathbf{E}|} \right), \tag{39}
 \end{aligned}$$

789 where  $\mu := 1 + 2c^{(2)}(1+\gamma)$  and we have removed the term  $-2c^{(2)}\gamma(1+\gamma)\sqrt{d_v d_u}/2|\mathbf{E}| \leq 0$ .  
 790 Moreover, since  $\lambda_1 < 1$  unless  $G$  is the complete graph (and if that was the case, then we could take  
 791 the distance  $r$  below to simply be equal to 1) and  $c_2 \leq 1$ , we find

$$\gamma \sqrt{d_v d_u} \frac{1}{2|\mathbf{E}|} \left( 1 - \frac{\mu}{c_2 \lambda_1} \left( 1 + \frac{|1 - c_2 \lambda^*|^{r-1}}{2} \left( \sqrt{\frac{d_v}{d_u}} + \sqrt{\frac{d_u}{d_v}} \right) \right) \right) \leq 0.$$

792 Accordingly, we can simplify (39) as

$$\begin{aligned} \text{mix}_{y_G^{(m)}}(v, u) &\leq \gamma \sqrt{d_v d_u} \left( \frac{m}{|E|} \mu - \frac{\tau(v, u)}{4c_2 |E|} \mu \right) + \frac{\gamma \mu}{2c_2 \lambda_1} \left( \sqrt{\frac{d_v}{d_u}} + \sqrt{\frac{d_u}{d_v}} \right) \\ &\quad + \frac{\gamma \mu}{c_2 \lambda_1} |1 - c_2 \lambda^*|^{r-1} + 2 \frac{\gamma c^{(2)}}{c_2}. \end{aligned}$$

793 We can now rearrange the terms and obtain

$$\begin{aligned} m &\geq \frac{\tau(v, u)}{4c_2} + \frac{|E|}{\sqrt{d_v d_u}} \left( \frac{\text{mix}_{y_G}(v, u)}{\gamma \mu} - \frac{1}{2c_2 \lambda_1} \left( \sqrt{\frac{d_v}{d_u}} + \sqrt{\frac{d_u}{d_v}} \right) - \frac{1}{c_2 \lambda_1} |1 - c_2 \lambda^*|^{r-1} - 2 \frac{c^{(2)}}{c_2} \frac{1}{\mu} \right) \\ &\geq \frac{\tau(v, u)}{4c_2} + \frac{|E|}{\sqrt{d_v d_u}} \left( \frac{\text{mix}_{y_G}(v, u)}{\gamma \mu} - \frac{1}{2c_2 \lambda_1} (2\gamma) - \frac{1}{c_2 \lambda_1} |1 - c_2 \lambda^*|^{r-1} - 2 \frac{c^{(2)}}{c_2} \frac{1}{\mu} \right) \\ &\geq \frac{\tau(v, u)}{4c_2} + \frac{|E|}{\sqrt{d_v d_u}} \left( \frac{\text{mix}_{y_G}(v, u)}{\gamma \mu} - \frac{1}{c_2} \left( \frac{\gamma + |1 - c_2 \lambda^*|^{r-1}}{\lambda_1} + 2 \frac{c^{(2)}}{\mu} \right) \right), \end{aligned}$$

794 which completes the proof.  $\square$

795 We note that the case of  $\mathbf{A} = \mathbf{A}_{\text{rw}}$  follows easily since one can adapt the previous argument exactly as  
796 in the proof of Theorem D.2, which lead to the same bounds once we replace  $\gamma$  with  $\gamma' = d_{\max}/d_{\min}$ .

797 First, we note that the bounds again simplify further and become sharper if the message-functions  $\psi$   
798 in (2) are linear.

799 **Corollary D.4.** *If the assumptions of Theorem 4.4 are satisfied, and the message-functions  $\psi$  are  
800 linear – as for GCN, GIN, GraphSAGE – then*

$$m \geq \frac{\tau(v, u)}{4c_2} + \frac{|E|}{\sqrt{d_v d_u}} \left( \frac{\text{mix}_{y_G}(v, u)}{\gamma} - \frac{1}{c_2 \lambda_1} \left( \gamma + |1 - c_2 \lambda^*|^{r-1} \right) \right).$$

#### 801 D.4 The case of the unnormalized adjacency matrix

802 In this Section we extend the analysis on the depth required to induce mixing, to the case of the  
803 unnormalized adjacency matrix  $\mathbf{A}$ . When  $\mathbf{A} = \mathbf{A}$ , the aggregation in (2) is simply a sum over the  
804 neighbours, a case that covers the classical GIN-architecture. In this way, the messages are no longer  
805 scaled down by the degree of (either) the endpoints of the edge, which means that, *in principle*, the  
806 whole GNN architecture is more sensitive but independent of where we are in the graph. First, we  
807 generalize Theorem 4.4 to this setting. We note that the same conclusions hold, provided that the  
808 maximal operator norm of the weights is smaller than the maximal degree  $d_{\max}$ ; this is not surprising,  
809 since it accounts for the lack of the normalization of the messages.

810 **Corollary D.5.** *Consider an MPNN as in (2) with  $\mathbf{A} = \mathbf{A}$ . If  $\omega/(wd_{\max}) + c_1 + c_2 \leq 1$  and  
811  $wd_{\max} \leq 1$ , then the minimal depth  $m$  satisfies the same lower bound as in Theorem 4.4 with  $\gamma = 1$ .*

812 *Proof.* First, we note that in this case

$$S_{ij} \leq \frac{\omega}{w} \delta_{ij} + c_1 d_{\max} \delta_{ij} + c_2 A_{ij} \leq d_{\max} \left( \alpha \delta_{ij} + c_2 (\mathbf{A}_{\text{sym}})_{ij} \right),$$

813 where  $\alpha = \omega/(wd_{\max}) + c_1$ . In particular, we find that

$$(\mathbf{S}^k \mathbf{1})_i \leq \sum_{p=0}^k \binom{k}{p} \alpha^{k-p} c_2^p (\mathbf{A}^p \mathbf{1})_i \leq (d_{\max})^k (\alpha + c_2)^k.$$

814 Accordingly, we can bound the first sum in (7) as

$$\left( \sum_{k=0}^{m-1} w^{2m-k} (d_{\max})^k (\alpha + c_2)^k (d_{\max})^{2(m-k)} \left( \alpha \mathbf{I} + c_2 \mathbf{A}_{\text{sym}} \right)^{2(m-k)} \right)_{vu} \leq \left( \sum_{k=0}^{m-1} \mathbf{Z}^{2(m-k)} \right)_{vu},$$

815 where we have used the assumptions  $\alpha + c_2 \leq 1$ , and  $wd_{\max} \leq 1$ , and the definition  $\mathbf{Z} := \mathbf{I} - c_2 \mathbf{\Delta}$   
816 in Theorem D.2. Since this term is the same one entering the argument in the proof of Theorem 4.4  
817 (once we set  $\gamma = 1$ ) we can proceed in the same way to estimate it. A similar argument works for the  
818 sum of the  $\mathbf{Q}_k$  terms, which, thanks to our assumptions, can still be bounded as in (30) with  $\gamma = 1$  so  
819 that we can finally simply copy the proof of Theorem 4.4.  $\square$

820 **A relative measurement for  $\widetilde{\text{OSQ}}$ .** To account for the fact that different message-passing matrices  
 821 **A** may lead to inherently quite distinct scales (think of the case where the aggregation is a mean vs  
 822 when it is a sum), one could modify the over-squashing characterization in Definition 4.2 as follows:

823 **Definition D.6.** Given an MPNN with capacity  $(m, w)$ , we define the **relative over-squashing of**  
 824 **nodes  $v, u$  as**

$$\widetilde{\text{OSQ}}_{v,u}^{\text{rel}}(m, w) := \left( \frac{\sum_{k=0}^{m-1} w^{2m-k-1} \left( w(\mathbf{S}^{m-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) \mathbf{S}^{m-k} + c^{(2)} \mathbf{Q}_k \right)_{vu}}{\max_{i,j \in \mathcal{V}} \sum_{k=0}^{m-1} w^{2m-k-1} \left( w(\mathbf{S}^{m-k})^\top \text{diag}(\mathbf{1}^\top \mathbf{S}^k) \mathbf{S}^{m-k} + c^{(2)} \mathbf{Q}_k \right)_{ij}} \right)^{-1}.$$

825 The normalization proposed here is similar to the idea of relative score introduced in [53]. This way,  
 826 a larger scale induced by a certain choice of the message-passing matrix **A**, is naturally accounted  
 827 for by the relative measurement. In particular, the relative over-squashing is now quantifying the  
 828 maximal mixing among a certain pair of nodes  $v, u$  compared to the maximal mixing that the same  
 829 MPNN over the same graph can generate among *any* pair of nodes. In our theoretical development  
 830 in Section 4 we have decided to rely on the absolute measurement since our analysis depends on  
 831 the derivation of the maximal mixing induced by an MPNN (i.e. upper bounds) which translate  
 832 into necessary criteria for an MPNN to generate a given level of mixing. In principle, to deal with  
 833 relative measurements, one would also need some form of lower bound on the maximal mixing and  
 834 hence address also whether the conditions provided are indeed sufficient. We reserve a thorough  
 835 investigation of this angle to future work.

## 836 **E The case of node-level tasks**

837 In this Section we discuss how one can extend our analysis to node-level tasks and further comment  
 838 on the novelty of our approach compared to existing results in [8, 21]. First, we emphasize that the  
 839 analysis on the Jacobian of node features carried over in [8, 21] cannot be extended to graph-level  
 840 functions and that in fact, *our notion of mixing is needed* to assess how two different node-features  
 841 are communicating when the target is a graph-level function.

842 From now on, let us consider the case where the function we need to learn is  $\mathbf{Y} : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ , and  
 843 as usual we assume it to be equivariant with respect to permutations of the nodes. A natural attempt  
 844 to connect the results in [8, 21] and the expressivity of the MPNNs – in the spirit of our Section 3 –  
 845 could be to characterize the **first-order interactions** (or mixing of order 1) of the features associated  
 846 with nodes  $v, u$  with respect to the underlying node-level task  $\mathbf{Y}$  as

$$\text{mix}_{\mathbf{Y}}^{(1)}(v, u) = \max_{\mathbf{X}} \max_{1 \leq \alpha, \beta \leq d} \left| \frac{\partial (\mathbf{Y}(\mathbf{X}))_v^\alpha}{\partial x_u^\beta} \right|,$$

847 where  $(\mathbf{Y}(\mathbf{X}))_v \in \mathbb{R}^d$  is the value of the node-level map at  $v$ . Accordingly, one can then use Theorem  
 848 C.1 to derive upper bounds on the maximal first-order interactions that MPNNs (2) can induce among  
 849 nodes. As a consequence of this approach, we would still find that MPNNs struggle to learn functions  
 850 with large  $\text{mix}_{\mathbf{Y}}^{(1)}(v, u)$  if nodes  $v, u$  have large commute time. In particular, in light of Theorem C.1,  
 851 we can extend the measure of over-squashing to the case of first-order interactions for node-level tasks.  
 852 Once again, below we tacitly assume that the non-linear activation  $\sigma$  satisfies  $|\sigma'| \leq 1$ , although it is  
 853 straightforward to extend the formulation to the general case.

854 **Definition E.1.** Given an MPNN as in (2) with capacity  $(m, w)$ , we define the **first-order over-**  
 855 **squashing of  $v, u$  as**

$$\text{OSQ}_{v,u}^{(1)}(m, w) := \left( \text{mix}_{\mathbf{Y}}^{(1)}(v, u) \right)^{-1}.$$

856 As for the case of graph-level tasks, we can then study a proxy (lower bound) for the node-level  
 857 over-squashing of order 1 by:

858 **Definition E.2.** Given an MPNN as in (2) with capacity  $(m, w)$  and **S** defined in (5), we approximate  
 859 the **first-order over-squashing of  $v, u$  as**

$$\widetilde{\text{OSQ}}_{v,u}^{(1)}(m, w) := \left( (c_\sigma w)^m (\mathbf{S}^m)_{vu} \right)^{-1}.$$

860 It follows then from Theorem C.1, that a necessary condition for an MPNN to learn a node-level  
 861 function  $\mathbf{Y}$  with first-order mixing  $\text{mix}_{\mathbf{Y}}^{(1)}(v, u)$  is

$$\widetilde{\text{OSQ}}_{v,u}^{(1)}(m, \mathbf{w}) < \left( \text{mix}_{\mathbf{Y}}^{(1)}(v, u) \right)^{-1}.$$

862 It is straightforward to argue as in Theorem 4.4 and [8] for example, to derive that nodes at higher  
 863 effective resistance will incur higher first-order over-squashing. Accordingly:

864 *An MPNN as in (2) with bounded capacity, cannot learn node-level functions with high first-order  
 865 interactions among nodes  $v, u$  with high effective resistance.*

866 **Building a hierarchy of measures.** Although first-order derivatives might be enough to capture  
 867 some form of over-squashing for node-level tasks, even in this scenario we can study the pairwise  
 868 mixing induced *at a specific node*, and hence consider the curvature (or Hessian) of the node-level  
 869 function  $\mathbf{Y}$  – which is more expressive than the first-order Jacobian. Accordingly, for a node-level  
 870 function  $\mathbf{Y} : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ , we say that it has **second-order interactions** (or mixing of order 2)  
 871  $\text{mix}_{\mathbf{Y}}^{(2)}(i, v, u)$  of the features associated with nodes  $v, u$  at a given node  $i$  when

$$\text{mix}_{\mathbf{Y}}^{(2)}(i, v, u) = \max_{\mathbf{X}} \left\| \frac{\partial^2 (\mathbf{Y}(\mathbf{X}))_i}{\partial \mathbf{x}_u \partial \mathbf{x}_v} \right\|.$$

872 We can then restate Theorem C.2 as follows – we let  $\mathbf{Y}^{(m)}$  be the node-level function computed by  
 873 an MPNN after  $m$  layers.

874 **Corollary E.3.** *Given MPNNs as in (2), let  $\sigma$  and  $\psi^{(t)}$  be  $\mathcal{C}^2$  functions and assume  $|\sigma'|, |\sigma''| \leq c_\sigma$ ,  
 875  $\|\Omega^{(t)}\| \leq \omega$ ,  $\|\mathbf{W}^{(t)}\| \leq \mathbf{w}$ ,  $\|\nabla_1 \psi^{(t)}\| \leq c_1$ ,  $\|\nabla_2 \psi^{(t)}\| \leq c_2$ ,  $\|\nabla^2 \psi^{(t)}\| \leq c^{(2)}$ . Let  $\mathbf{S} \in \mathbb{R}^{n \times n}$  be  
 876 defined as in (5). Given nodes  $i, v, u \in \mathbf{V}$ , if  $\mathbf{P}_{(vu)}^{(\ell)} \in \mathbb{R}^n$  is as in (14) and  $m$  is the number of layers,  
 877 then the maximal mixing of order 2 of the MPNN at node  $i$  satisfies*

$$\begin{aligned} \text{mix}_{\mathbf{Y}^{(m)}}^{(2)}(i, v, u) &\leq \sum_{k=0}^{m-1} \sum_{j \in \mathbf{V}} (c_\sigma \mathbf{w})^{2m-k-1} \mathbf{w} (\mathbf{S}^{m-k})_{jv} (\mathbf{S}^k)_{ij} (\mathbf{S}^{m-k})_{ju} \\ &\quad + c^{(2)} \sum_{\ell=0}^{m-1} (c_\sigma \mathbf{w})^{m+\ell} (\mathbf{S}^{m-1-\ell} \mathbf{P}_{(vu)}^{(\ell)})_i. \end{aligned} \quad (40)$$

878 Similarly to Definition 4.2, we can use the maximal mixing (at the node-level) to characterize the  
 879 over-squashing of order two at a specific node as follows: as usual, for simplicity we assume that  
 880  $c_\sigma = 1$ .

881 **Definition E.4.** *Given an MPNN as in (2) with capacity  $(m, \mathbf{w})$  and  $\mathbf{S}$  defined in (5), we approximate  
 882 the **second-order over-squashing** of  $v, u$  at node  $i$  as*

$$\begin{aligned} \widetilde{\text{OSQ}}_{i,v,u}^{(2)}(m, \mathbf{w}) &:= \left( \sum_{k=0}^{m-1} \sum_{j \in \mathbf{V}} \mathbf{w}^{2m-k} (\mathbf{S}^{m-k})_{jv} (\mathbf{S}^k)_{ij} (\mathbf{S}^{m-k})_{ju} \right. \\ &\quad \left. + c^{(2)} \sum_{\ell=0}^{m-1} \mathbf{w}^{m+\ell} (\mathbf{S}^{m-1-\ell} \mathbf{P}_{(vu)}^{(\ell)})_i \right)^{-1}. \end{aligned}$$

883 It is then straightforward to extend our theoretical analysis to derive how  $\widetilde{\text{OSQ}}^{(2)}$  prevents MPNNs  
 884 from learning node-level functions with high-mixing at some specific node  $i$  of features associated  
 885 with nodes  $v, u$  at large commute time. To support our claim, consider the setting in Theorem 4.4 and  
 886 hence let  $\mathbf{A} = \mathbf{A}_{\text{sym}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ . Under the same assumptions of Theorem 4.4, we find

$$(\mathbf{S}^k)_{ij} \leq 1.$$

887 We can then simply copy the proof of Theorem 4.4 once we set  $\gamma = 1$  and extend its conclusions as  
 888 follows:

889 **Corollary E.5.** Consider an MPNN as in (2) and let the assumptions of Theorem 4.4 hold. If the  
 890 MPNN generates second-order mixing  $\text{mix}_{\mathcal{V}}^{(2)}(i, v, u)$  at node  $i$ , with respect to the features associated  
 891 with nodes  $v, u$ , then the number of layers  $m$  satisfy:

$$m \geq \frac{\tau(v, u)}{4c_2} + \frac{|\mathcal{E}|}{\sqrt{d_v d_u}} \left( \frac{\text{mix}_{y_G}(v, u)}{\mu} - \frac{1}{c_2} \left( \frac{1 + |1 - c_2 \lambda^*|^{r-1}}{\lambda_1} + 2 \frac{c^{(2)}}{\mu} \right) \right),$$

892 where  $\mu = 1 + 4c^{(2)}$ .

893 Accordingly, in this Section we have adapted our results from graph-level tasks to node-level tasks  
 894 and proved that:

895 **The message of Section E.** An MPNN of bounded capacity  $(m, w)$ , cannot learn node-level functions  
 896 that, at some node  $i$ , induce high (first order or second order) mixing of features associated with  
 897 nodes  $v, u$  whose commute time is large.

## 898 F Additional details of experiments and further ablations 5

### 899 F.1 The role of mixing

900 We further test the considered MPNN architectures on their performance with respect to different  
 901 mixings. To this end, we consider again the tanh-based mixing as in our previous tasks (i.e., regressing  
 902 targets  $y_i = \tanh(x_{u^i}^i + x_{v^i}^i)$  for each graph  $G^i$  in the dataset), as well as another mixing based on  
 903 the exponential function (i.e., with targets  $y_i = \exp(x_{u^i}^i + x_{v^i}^i)$ ). We note that these two tasks differ  
 904 significantly in terms of their maximal mixing values (4) (shown in Table 1). Thus, according to (8)  
 905 and Theorem 4.4, we would expect that MPNNs perform significantly worse in the case of higher  
 906 maximal mixing, i.e., for the exponential-based mixing compared to the tanh-mixing. To confirm  
 907 this empirically, we train the MPNNs on both types of mixing and provide the resulting relative  
 908 MAEs (i.e., MAE divided by the  $L^1$ -norm of the targets) in Table 1. We can see that all four MPNNs  
 909 perform significantly better on the tanh-mixing than on the exponential-based mixing. Moreover,  
 910 increasing the range for the exponential-based mixing from 1 to 1.5 further impairs the performance  
 911 of all considered MPNNs. In order to check if this difference in performance can simply be explained  
 912 by a higher capacity required by a neural network to accurately approximate the mapping  $\exp(x + y)$   
 913 compared to  $\tanh(x + y)$  for some inputs  $x, y \in \mathbb{R}$ , we train a simple two-layer feed-forward neural  
 914 network (with 2 inputs, i.e.,  $x$  and  $y$ ) on both mappings. The trained networks reach a similarly low  
 915 relative MAE of  $4.6 \times 10^{-4}$  for the  $\tanh(x + y)$  mapping as well as  $4.1 \times 10^{-4}$  for the  $\exp(x + y)$   
 916 mapping using an input range of  $(0, 1)$  and  $4.0 \times 10^{-4}$  for an input range of  $(0, 1.5)$ . Thus, we can  
 917 conclude that the significant differences in the obtained results in Table 1 are not caused by a higher  
 918 capacity required by a neural network to learn the underlying mappings of the different mixings.

Table 1: Relative MAE of GCN, GIN, GraphSAGE and GatedGCN on different choices of mixing on synthetic ZINC for a fixed 0.8-quantile of the commute time distributions over graphs  $G_i$ .

| Mixing                         | input interval | maximal mixing | GCN   | GIN   | GraphSAGE | GatedGCN |
|--------------------------------|----------------|----------------|-------|-------|-----------|----------|
| $\tanh(x_{u^i}^i + x_{v^i}^i)$ | $(0, 1)$       | $\approx 0.77$ | 0.024 | 0.014 | 0.006     | 0.004    |
| $\exp(x_{u^i}^i + x_{v^i}^i)$  | $(0, 1)$       | $\approx 7.4$  | 0.043 | 0.021 | 0.033     | 0.008    |
| $\exp(x_{u^i}^i + x_{v^i}^i)$  | $(0, 1.5)$     | $\approx 20.1$ | 0.054 | 0.035 | 0.075     | 0.014    |

### 919 F.2 Computing the commute time

920 The commute time  $\tau$  between two nodes  $u, v \in \mathcal{V}$  on a graph  $G$  can be efficiently computed via the  
 921 effective resistance  $R$ , with  $\tau(u, v) = 2|\mathcal{E}|R(u, v)$ . In order to compute the effective resistance  $R$ , we  
 922 introduce the (non-normalized) Laplacian matrix  $\mathbf{L} = \mathbf{D} - \mathbf{A}$ , where  $\mathbf{D}$  is the degree matrix. The  
 923 effective resistance can then be computed by

$$R(u, v) = \Gamma_{uu} + \Gamma_{vv} - 2\Gamma_{uv},$$

924 where  $\mathbf{\Gamma}$  is the the Moore-Penrose inverse of

$$\mathbf{L} + \frac{1}{|\mathcal{V}|} \mathbf{1}_{|\mathcal{V}| \times |\mathcal{V}|},$$



925 with  $\mathbf{1}_{|V|\times|V|} \in \mathbb{R}^{|V|\times|V|}$  being a matrix with all entries set to one.

926 **F.3 On the training error**

927 In this section, we report the training error of the MPNNs trained in section 5. Fig. 5 shows the  
 928 training MAE corresponding to the experiment in section 5.1, while Fig. 6 shows the training MAE  
 929 corresponding to section 5.2. We can see that in both cases, the training MAE exhibits the same  
 930 qualitative behavior as the reported test MAE in the main paper, i.e., the training MAE increases  
 931 for increasing levels of commute time  $\tau$ , while it decreases for increasing number of MPNN layers,  
 which further validates our claim that *over-squashing hinders the expressive power of MPNNs*.

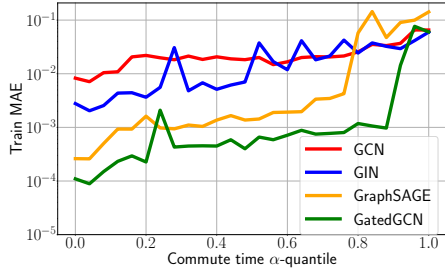


Figure 5: Train MAE of GCN, GIN, GraphSAGE, and GatedGCN on synthetic ZINC, where the commute time of the underlying mixing is varied, while the MPNN architecture is fixed (e.g., depth, number of parameters), i.e., mixing according to increasing values of the  $\alpha$ -quantile of the  $\tau$ -distribution over the ZINC graphs.

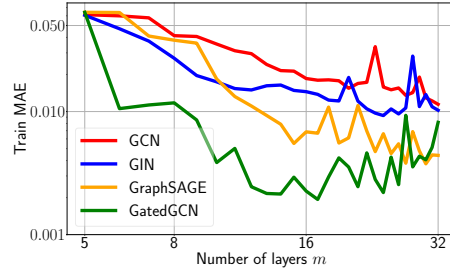


Figure 6: Train MAE of GCN, GIN, GraphSAGE, and GatedGCN on synthetic ZINC, where the commute time is fixed to be high (i.e., at the level of the 0.8-quantile), while only the depth of the underlying MPNN is varied between 5 and 32 (all other architectural components are fixed).

932

933 **F.4 On the over-squashing measure**

934 In this section, we examine how the over-squashing measure  $\widetilde{\text{OSQ}}$  (as of Definition 4.2) depends on  
 935 the commute time  $\tau$  as well as on the depth of the underlying MPNN. To this end, we follow the  
 936 experimental setup of section 5.1 and 5.2, but instead of training the models and presenting their  
 937 performance in terms of the test MAE, we compute  $\widetilde{\text{OSQ}}$  of the underlying models. We can see in  
 938 Fig. 7 that  $\widetilde{\text{OSQ}}$  increases for increasing values of the  $\alpha$ -quantile of the  $\tau$ -distribution for all MPNNs  
 939 considered here. Moreover, we can see in Fig. 8 that  $\widetilde{\text{OSQ}}$  decreases for increasing number of layers  
 940 for all considered models.

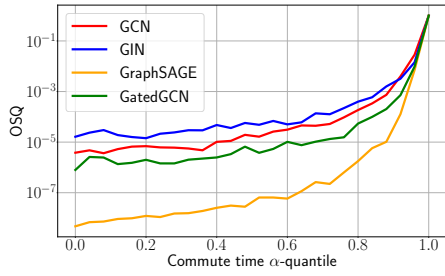


Figure 7:  $\widetilde{\text{OSQ}}$  (Definition 4.2) of GCN, GIN, GraphSAGE, and GatedGCN on synthetic ZINC, where the commute time of the underlying mixing is varied, while the MPNN architecture is fixed (e.g., depth, number of parameters), i.e., mixing according to increasing values of the  $\alpha$ -quantile of the  $\tau$ -distribution over the ZINC graphs.

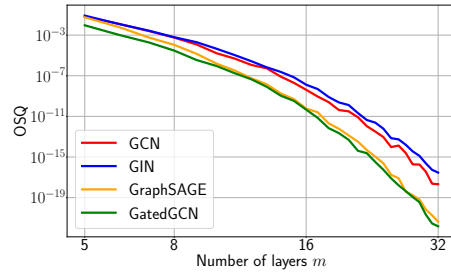


Figure 8:  $\widetilde{\text{OSQ}}$  (Definition 4.2) of GCN, GIN, GraphSAGE, and GatedGCN on synthetic ZINC, where the commute time is fixed to be high (i.e., at the level of the 0.8-quantile), while only the depth of the underlying MPNN is varied between 5 and 32 (all other architectural components are fixed).