

# BREAKING THE EXPRESSIVE BOTTLENECKS OF GRAPH NEURAL NETWORKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recently, the Weisfeiler-Lehman (WL) graph isomorphism test was used to measure the expressiveness of graph neural networks (GNNs), showing that the neighborhood aggregation GNNs were at most as powerful as 1-WL test in distinguishing graph structures. There were also improvements proposed in analogy to  $k$ -WL test ( $k > 1$ ). However, the aggregators in these GNNs are far from injective as required by the WL test, and suffer from weak distinguishing strength, making it become expressive bottlenecks. In this paper, we improve the expressiveness by exploring powerful aggregators. We reformulate aggregation with the corresponding aggregation coefficient matrix, and then systematically analyze the requirements of the aggregation coefficient matrix for building more powerful aggregators and even injective aggregators. It can also be viewed as the strategy for preserving the rank of hidden features, and implies that basic aggregators correspond to a special case of low-rank transformations. We also show the necessity of applying non-linear units ahead of aggregation, which is different from most aggregation-based GNNs. Based on our theoretical analysis, we develop two GNN layers, ExpandingConv and CombConv. Experimental results show that our models significantly boost performance, especially for large and densely connected graphs.

## 1 INTRODUCTION

Graphs are ubiquitous in the real world. Social networks, traffic networks, knowledge graphs, and molecular structures are typical graph-structured data. Graph Neural Networks (GNNs) (Scarselli et al., 2008; Gori et al., 2005), leveraging the power of neural networks to graph-structured data, have a rapid development recently (Kipf & Welling, 2016; Hamilton et al., 2017; Bronstein et al., 2017; Gilmer et al., 2017; Duvenaud et al., 2015).

Expressive power of GNNs measures their abilities to represent different graph structures (Sato, 2020). It decides the performance of GNNs where the awareness of graph structures is required, especially on large graphs with complex topologies. The neighborhood aggregation scheme (or message passing) follows the same pattern with weisfeiler-lehman (WL) graph isomorphism test (Weisfeiler & Leman, 1968) to encode graph structures, where node representations are computed iteratively by aggregating transformed representations of its neighbors with structural information learned implicitly. Therefore, the WL test is used to measure the expressiveness of GNNs. Unfortunately, general GNNs are at most as powerful as 1-order WL test (Morris et al., 2019; Xu et al., 2019). There is also work trying to improve the expressiveness that are beyond 1-order WL test (Maron et al., 2019; Morris et al., 2019; Chen et al., 2019; Li et al., 2020b; Vignac et al., 2020). However, the weak distinguishing strength of aggregators is the fundamental limitation. The expressiveness analysis measured by the WL test assumes that aggregators are injective, which is usually unreachable. Therefore, this motivates us to investigate the following questions: What are the key factors to limit the expressiveness of GNN? and how to break these limitations?

Aggregators are permutation invariant functions that operate on sets while preserving permutation invariance. (Zaheer et al., 2017) first theoretically studied permutation invariant functions and provided a family of functions to which any permutation invariant function must belong. (Xu et al., 2019) extended it on multisets but only for countable space. (Corso et al., 2020) further extended it to uncountable space. (Murphy et al., 2018) and (Murphy et al., 2019) expressed a permutation invariant function by approximating an average over permutation-sensitive functions with tractability

strategies. (Dehmamy et al., 2019) showed that a single propagation rule applied in general GNNs is rather restrictive in learning graph moments (Lin & Skiena, 1995). They and (Corso et al., 2020) improved the distinguishing strength of aggregation by leveraging multiple basic aggregators (SUM, MEAN, NORMALIZED MEAN, MAX/MIN, and STD). This strategy showed its effectiveness on tasks taken from classical graph theory.

In contrast to existing studies towards aggregators in GNNs, we provide a new GNN formulation, where the aggregation is represented as the multiplication of the corresponding hidden feature matrix of neighbors and the aggregation coefficient matrix. This new formulation enables us to answer the following questions: (i) when a GNN will lose its expressive power; (ii) How to build aggregators with higher distinguishing strength, even injective aggregators. Based on our theoretical analysis, we propose two GNN layers: ExpandingConv and CombConv, and evaluate them on general graph classification and graph regression tasks. Our key contributions are summarized as follows:

- We formalize the distinguishing strength of aggregators as a partial order, and theoretically show that the choice of aggregators can be bottlenecks of expressiveness. We also propose to apply nonlinear units ahead of aggregations to break the distinguishing strength limitations of aggregators as well as to achieve an implicit sampling mechanism.
- We reformulate the neighborhood aggregation with the aggregation coefficient matrix and then provide a theoretical point of view on building powerful aggregators and even injective aggregators.
- We propose ExpandingConv and CombConv layers which achieve state-of-the-art performance on a variety of graph tasks. We also show that multi-head GAT is one of the ExpandingConv implementations, which brings a theoretical explanation for its effectiveness.

## 2 PRELIMINARIES

### 2.1 NOTATIONS

For a graph  $G(V, E)$ , we denote the set of edges, nodes and node feature vectors respectively by  $E_G$ ,  $V_G$  and  $X_G$ .  $\mathcal{N}(v)$  represents the set of neighbors of  $v$  including itself, i.e.,  $\mathcal{N}(v) = \{u \in V_G | (u, v) \in E_G\} \cup \{v\}$ . We use  $[n]$  to denote the set  $\{1, 2, \dots, n\}$ .  $\{\{\dots\}\}$  represents a multi-set, i.e., a set with possibly repeating elements.  $\Pi_n$  represents the set of all permutations of the integers 1 to  $n$ .  $\mathbf{h}_\pi$ , where  $\pi \in \Pi_{|h|}$ , is a reordering of the elements of a sequence  $\mathbf{h}$  according to  $\pi$ . Given a matrix  $\mathbf{X} \in \mathbb{R}^{a \times b}$ ,  $\mathbf{X}^T$  represents the transpose of  $\mathbf{X}$ , and  $\text{vec}(\mathbf{X}) \in \mathbb{R}^{ab \times 1}$  represents the column stack of  $\mathbf{X}$ .

### 2.2 GRAPH NEURAL NETWORKS

Most GNNs adopt the neighborhood aggregation scheme (Gilmer et al., 2017) to learn the node representations, which utilizes both node features and graph structures. In the  $k$ -th layer, the representation of node  $v$

$$\mathbf{h}_v^{(k)} = \text{Update}(\mathbf{h}_v^{(k-1)}, \text{Aggregate}(\{\{\mathbf{h}_u^{(k-1)} | u \in \mathcal{N}(v)\}\})). \quad (1)$$

**Aggregators in GNNs.** An aggregator is a permutation invariant function (Zaheer et al., 2017) with bounded size inputs. It satisfies: (i) size insensitive: an aggregator can take an arbitrary but finite size of inputs; (ii) permutation invariant: an aggregator is invariant to the permutation of input. There are a limited number of basic aggregators such as SUM, MEAN, NORMALIZED MEAN, MAX/MIN, STD, etc. Most proposed GNNs apply one of these aggregators. Sum-of-power mapping (Zaheer et al., 2017) and normalized moments (Corso et al., 2020) can also be used as aggregators and they allow for a variable number of aggregators.

## 3 PROPOSED MODEL

In this section, we first formalize the distinguishing strength of aggregators as a partial order, and show why basic aggregators used in popular GNNs become bottlenecks of expressiveness. Then, we analyze the requirements for building powerful aggregators and even injective aggregators. Finally, we introduce two GNN layers based on our theoretical analysis.

### 3.1 DISTINGUISHING STRENGTH OF AGGREGATORS

To ensure generality, our analysis of aggregators is always considered in multisets and uncountable case, where the inputs are continuous and with possibly repeating elements. We first introduce distinguishing strength under the concept of partial order (Schmidt, 2011).

**Distinguishing strength.** The distinguishing strength of aggregator  $f_{aggr1}$  is stronger than  $f_{aggr2}$ , denoted by  $f_{aggr1} \geq f_{aggr2}$ , if and only if for any two multisets  $\mathbf{x}_1$  and  $\mathbf{x}_2$  where the number of elements can be different,  $f_{aggr2}(\mathbf{x}_1) \neq f_{aggr2}(\mathbf{x}_2) \Rightarrow f_{aggr1}(\mathbf{x}_1) \neq f_{aggr1}(\mathbf{x}_2)$ . Meanwhile, if there exist  $\mathbf{x}'_1$  and  $\mathbf{x}'_2$  such that  $f_{aggr1}(\mathbf{x}'_1) \neq f_{aggr1}(\mathbf{x}'_2)$  but  $f_{aggr2}(\mathbf{x}'_1) = f_{aggr2}(\mathbf{x}'_2)$ ,  $f_{aggr1}$  is strictly stronger than  $f_{aggr2}$ , denoted by  $f_{aggr1} > f_{aggr2}$ . If  $f_{aggr1} \geq f_{aggr2} \geq f_{aggr1}$ , we say these two aggregators have the same distinguishing strength, denoted by  $f_{aggr1} = f_{aggr2}$ . If there exist multisets  $\mathbf{x}_1$  and  $\mathbf{x}_2$  such that  $f_{aggr1}(\mathbf{x}_1) \neq f_{aggr1}(\mathbf{x}_2)$ ,  $f_{aggr2}(\mathbf{x}_1) = f_{aggr2}(\mathbf{x}_2)$ , and there also exist  $\mathbf{x}'_1$  and  $\mathbf{x}'_2$  such that  $f_{aggr1}(\mathbf{x}'_1) = f_{aggr1}(\mathbf{x}'_2)$ ,  $f_{aggr2}(\mathbf{x}'_1) \neq f_{aggr2}(\mathbf{x}'_2)$ , we say  $f_{aggr1}$  and  $f_{aggr2}$  are incomparable.

Distinguishing strength is a partial order, and the set of all aggregators form a poset. In this poset, the aggregators with the greatest distinguishing strength should be injective. With the definition of distinguishing strength, we can compare any two aggregators. The distinguishing strength of widely used aggregators SUM, MEAN, MAX/MIN is incomparable. One can easily give two multisets of elements that are distinguished by one aggregator but are not distinguished by the others as showed in (Corso et al., 2020).

**Equivariant aggregator.**  $f_{aggr} : \{\{\mathbb{R}^d\}\} \rightarrow \mathbb{R}^d$  is an equivariant aggregator if and only if  $f_{aggr}(\{\{T \cdot \mathbf{x}_1, T \cdot \mathbf{x}_2, \dots, T \cdot \mathbf{x}_n\}\}) = T \cdot f_{aggr}(\{\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}\})$  for any  $T \in \mathbb{R}^{m \times d}$  and  $\{\{\mathbf{x}_i \in \mathbb{R}^d | i \in [n]\}\}$ .

Widely used SUM and MEAN are equivariant aggregators but MAX/MIN is not. We denote  $f_{aggr1} \otimes f_{aggr2}$  a new aggregator by combing  $f_{aggr1}$  and  $f_{aggr2}$  with  $f_{aggr1} \otimes f_{aggr2}(X) = [f_{aggr1}(X) || f_{aggr2}(X)]$ , where  $||$  denotes concatenation.

**Lemma 1.** (i) For any continuous function  $g$ , we have  $g \circ f_{aggr} \leq f_{aggr}$ , and when  $g$  is injective,  $f_{aggr} = g \circ f_{aggr}$ ; (ii)  $f_{aggr1} \otimes f_{aggr2} \geq f_{aggr1}$  and  $f_{aggr1} \otimes f_{aggr2} \geq f_{aggr2}$ . If  $f_{aggr1}$  and  $f_{aggr2}$  are incomparable,  $f_{aggr1} \otimes f_{aggr2} > f_{aggr1}$  and  $f_{aggr1} \otimes f_{aggr2} > f_{aggr2}$ ; (iii) If  $f_{aggr}$  is an equivariant aggregator, then  $f_{aggr}(T \cdot \mathbf{x}_1, T \cdot \mathbf{x}_2, \dots, T \cdot \mathbf{x}_n) \leq f_{aggr}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  for any  $T \in \mathbb{R}^{m \times d}$  and  $\{\{\mathbf{x}_i \in \mathbb{R}^d | i \in [n]\}\}$ .

We prove Lemma 1 in Appendix B. Lemma 1 indicates that aggregators become bottlenecks of distinguishing strength. For the equivariant aggregator, any linear transformation before aggregation and any transformation after aggregation have no contribution to the distinguishing strength. For SUM and MEAN, we have  $g(\text{SUM}(T \cdot \mathbf{x}_1, T \cdot \mathbf{x}_2, \dots, T \cdot \mathbf{x}_n)) \leq \text{SUM}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  and  $g(\text{MEAN}(T \cdot \mathbf{x}_1, T \cdot \mathbf{x}_2, \dots, T \cdot \mathbf{x}_n)) \leq \text{MEAN}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ , where  $T \in \mathbb{R}^{m \times d}$ , and  $g$  can be any continuous function. Based on Lemma 1, we can now compare the distinguishing strength of aggregations in some popular GNNs. GIN-0 sums all hidden features of neighbors at first, and then pass them to a 2-layer MLP. Therefore, when considering in a continuous input features space, the distinguishing strength of GIN-0 is at most as powerful as the SUM aggregator. GCN uses a NORMALIZED MEAN (denoted by  $n\text{MEAN}$ ) aggregator. Given a node  $v$  and its neighbors,  $n\text{MEAN}(v, u_1, \dots, u_n) = \frac{1}{\sqrt{|\mathcal{N}(v)|}} \cdot (\frac{h_v}{\sqrt{|\mathcal{N}(v)|}} + \frac{h_{u_1}}{\sqrt{|\mathcal{N}(u_1)|}} + \dots + \frac{h_{u_n}}{\sqrt{|\mathcal{N}(u_n)|-1}})$ .  $n\text{MEAN}$  is also an equivariant aggregator, and the distinguishing strength of aggregation in GCN is at most as powerful as  $n\text{MEAN}$ . GAT corresponds to the weighted SUM aggregation, where the weight coefficients are the functions of hidden features. This makes the distinguishing strength of GAT and SUM incomparable. Based on these observations, a potential approach to breaking the distinguishing strength limitation is to apply a nonlinear processing on inputs before aggregation.

### 3.2 BUILDING POWERFUL AGGREGATORS

In this section, we analyze the requirements for building more powerful aggregators and further injective aggregators. We first introduce a new representation of GNN layers which unifies several popular GNN layers. Given a node  $v$  and its neighbors  $\mathcal{N}(v)$ , our new formulation represents the

GNN operation as follows:

$$\begin{cases} \mathbf{m}_v = \mathbf{f}_{\text{local}}(v) & /* \text{aggregation coefficients generation} */ \\ \mathbf{r}_v^{(t)} = \mathbf{m}_{v\pi} \bar{\mathbf{h}}_{v\pi}^{(t-1)T} & /* \text{neighborhood aggregation} */ \\ \mathbf{h}_v^{(t)} = \mathbf{f}_{\text{NN}}(\mathbf{r}_v^{(t)}) & /* \text{feature/structure extraction} */ \end{cases} \quad (2)$$

Here,  $\mathbf{m}_v \in \mathbb{R}^{|\mathcal{N}(v)|}$  is the aggregation coefficients vector of node  $v$ . Note that  $\mathbf{m}_v$  should be the mapping of local structures such as node degrees, node or edge features of the  $k$ -hop neighbors assigned on node  $v$  to ensure the same encoding of isomorphic graphs.  $\bar{\mathbf{h}}_{v\pi}^{(t-1)T} = (\mathbf{h}_v^{(t-1)}, \mathbf{h}_{u_1}^{(t-1)}, \dots, \mathbf{h}_{u_{|\mathcal{N}(v)|}}^{(t-1)})^T \in \mathbb{R}^{|\mathcal{N}(v)| \times d}$  is the matrix representation of  $v$ 's neighbors according to a permutation  $\pi$ .  $\mathbf{f}_{\text{NN}}: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  is a neural network that extracts task-relevant information from the aggregated representation  $\mathbf{r}_v^{(t)}$ , and is used to update hidden feature  $\mathbf{h}_v^{(t)}$  of node  $v$ .

According to Equation 2, the aggregation should be with high distinguishing strength to avoid indistinguishability among neighbors. Meanwhile, the extraction should be powerful enough to efficiently extract task-relevant structural patterns from the aggregated representation of neighbors. Based on these observations, we reformulate GCN, GIN0 and GAT with their corresponding three-stage representations as follows:

$$\begin{aligned} \text{GCN} : \begin{cases} \mathbf{m}_v = \frac{1}{\sqrt{|\mathcal{N}(v)|}} \left( \frac{1}{\sqrt{|\mathcal{N}(v)|}}, \frac{1}{\sqrt{|\mathcal{N}(u_1)|}}, \dots, \frac{1}{\sqrt{|\mathcal{N}(u_{|\mathcal{N}(v)|-1})|}} \right) \\ \mathbf{r}_v^{(t)} = \mathbf{m}_{v\pi} \bar{\mathbf{h}}_{v\pi}^{(t-1)T} \\ \mathbf{h}_v^{(t)} = \sigma(\mathbf{W} \mathbf{r}_v^{(t)} + \mathbf{b}), \end{cases} & \text{GIN0} : \begin{cases} \mathbf{m}_v = \mathbf{1}^{1 \times |\mathcal{N}(v)|} \\ \mathbf{r}_v^{(t)} = \mathbf{m}_{v\pi} \bar{\mathbf{h}}_{v\pi}^{(t-1)T} \\ \mathbf{h}_v^{(t)} = \text{MLP}(\mathbf{r}_v^{(t)}), \end{cases} \\ \text{GAT} : \begin{cases} \mathbf{m}_v^{(t)} = (\text{att}(\mathbf{h}_v^{(t-1)}, \mathbf{h}_v^{(t-1)}), \text{att}(\mathbf{h}_v^{(t-1)}, \mathbf{h}_{u_1}^{(t-1)}), \dots, \text{att}(\mathbf{h}_v^{(t-1)}, \mathbf{h}_{u_{|\mathcal{N}(v)-1}}^{(t-1)})) \\ \mathbf{r}_v^{(t)} = \mathbf{m}_{v\pi}^{(t)} \bar{\mathbf{h}}_{v\pi}^{(t-1)T} \\ \mathbf{h}_v^{(t)} = \sigma(\mathbf{W} \mathbf{r}_v^{(t)} + \mathbf{b}). \end{cases} \end{aligned}$$

Their default formulations are given in Appendix A. In the aggregation step, GCN's  $\mathbf{m}_v$  is the mapping of neighbors' degrees; GIN0's  $\mathbf{m}_v$  is the mapping of node  $v$ 's degree which is equivalent to SUM aggregator; GAT's  $\mathbf{m}_v$  is the mapping of neighbors' features. All of them are the mappings of local structures as given in Equation 2.

In this three-stage representation, the aggregation is reformulated as the multiplication of the aggregation coefficients vector and the feature matrix of neighbors. It provides insights on improving the distinguishing strength of aggregations. First, we show how to characterize the permutation invariance in this formulation. Let  $\mathbf{M} \in \mathbb{R}^{s \times n}$  denote an aggregation coefficient matrix where  $s \geq 1$ . Note that in GCN, GIN and GAT,  $s$  is restricted to be 1.  $\mathbf{h}_\pi \in \mathbb{R}^{n \times d}$  is the matrix representation of  $n$  input elements according to  $\pi$ . The aggregation computation in the second step of Equation 2 is  $\mathbf{f}_{\text{aggr}}(\mathbf{M}, \mathbf{h}_\pi) = \mathbf{M} \mathbf{P}_\pi \mathbf{h}_\pi = \mathbf{M}_\pi \mathbf{h}_\pi$ , where  $\mathbf{P}_\pi$  is the permutation matrix according to  $\pi$ .  $\mathbf{P}_\pi \mathbf{h}_\pi$  ensures the same output for all  $\mathbf{h}_\pi$ ,  $\pi \in \Pi_{|h|}$ .  $\mathbf{M}_\pi = \mathbf{M} \mathbf{P}_\pi$  is the reordering of columns of  $\mathbf{M}$  according to  $\pi$ . For any  $\pi_1, \pi_2 \in \Pi_n$ ,  $\mathbf{f}_{\text{aggr}}(\mathbf{M}, \mathbf{h}_{\pi_1}) = \mathbf{f}_{\text{aggr}}(\mathbf{M}, \mathbf{h}_{\pi_2})$ , thus permutation invariance holds. Once  $\mathbf{M}$  is decided, we obtain a unique aggregator denoted by  $\mathbf{f}_\mathbf{M}$ . For any sequence of input elements  $\mathbf{h}$ ,  $\mathbf{f}_\mathbf{M}(\mathbf{h}) = \mathbf{f}_{\text{aggr}}(\mathbf{M}, \mathbf{h}_\pi)$ , where  $\pi \in \Pi_n$  can be any ordering of neighbors. Next, we analyze the distinguishing strength of  $\mathbf{f}_\mathbf{M}$ .

**Proposition 1.** *For any two matrices  $\mathbf{M} \in \mathbb{R}^{s \times n}$  and  $\mathbf{M}' \in \mathbb{R}^{s' \times n}$  with  $s, s' \leq n$ , we have (i)  $\mathbf{f}_{\begin{pmatrix} \mathbf{M} \\ \mathbf{M}' \end{pmatrix}} \geq \mathbf{f}_\mathbf{M}$ , where  $\begin{pmatrix} \mathbf{M} \\ \mathbf{M}' \end{pmatrix}$  means stacking these two matrices; (ii)  $\mathbf{f}_{\begin{pmatrix} \mathbf{M} \\ \mathbf{M}' \end{pmatrix}} > \mathbf{f}_\mathbf{M}$  if and only if  $\text{rank}(\begin{pmatrix} \mathbf{M} \\ \mathbf{M}' \end{pmatrix}) > \text{rank}(\mathbf{M})$ ; (iii) Any multiset of size  $n$  is distinguishable with  $\mathbf{f}_\mathbf{M}$  if and only if  $\text{rank}(\mathbf{M}) = n$ .*

We prove Proposition 1 in Appendix C. Proposition 1 shows that the distinguishing strength of  $\mathbf{f}_\mathbf{M}$  is decided by the rank of the corresponding  $\mathbf{M}$ . Yet, the distinguishing strength analysis in Proposition 1 is only suitable for multisets aggregated with shared  $\mathbf{f}_\mathbf{M}$ . Next, we extend the analysis for the case of different aggregators.

Let  $\text{Res}(\mathbf{f}_\mathbf{M})$  denote the set of all outputs of  $\mathbf{f}_\mathbf{M}$ . Our proposed three-stage representation also provides useful insight on the constraints among different aggregators. That is, in order to fully distinguish different local structures, for any two different  $\mathbf{f}_{\mathbf{M}_1}$  and  $\mathbf{f}_{\mathbf{M}_2}$ ,  $\text{Res}(\mathbf{f}_{\mathbf{M}_1}) \cap \text{Res}(\mathbf{f}_{\mathbf{M}_2}) = \emptyset$ . This is because to fully distinguish different local structures, we should ensure their aggregated

representations are different. Since  $M$  is restricted to be the mapping of local structures (such as  $k$ -hop neighbors), different  $M$  means that the corresponding local structures are different. Therefore, the aggregation results of different  $f_M$  must be different. However, it is not satisfied by existing GNNs, and there are few studies on distinguishing multisets aggregated by different aggregators. In Proposition 2, we present a detailed analysis of it.

**Proposition 2.** For any  $M_1, M_2 \in \mathbb{R}^{s \times n_1}$  and  $M'_1, M'_2 \in \mathbb{R}^{s' \times n_2}$ , (i)  $Res(f_{\binom{M_1}{M'_1}}) \cap Res(f_{\binom{M_2}{M'_2}}) \subseteq Res(f_{M_1}) \cap Res(f_{M_2})$ ; (ii) If  $Res(f_{\binom{M_1}{M'_1}}) \cap Res(f_{\binom{M_2}{M'_2}}) \subset Res(f_{M_1}) \cap Res(f_{M_2})$ , then  $rank(\binom{M_1 \ M_2}{M'_1 \ M'_2}) > rank(\binom{M_1 \ M_2}{M_1 \ M_2})$ ;

We prove Proposition 2 in Appendix D. Proposition 2 shows the necessity of preserving the rank of aggregation coefficient matrix when considering the distinguishing strength among different aggregators. Next, we provide a sufficient condition for building desired multiple injective aggregators with the outputs having no intersections.

**Proposition 3.** For any two aggregators  $f_{M_1}$  and  $f_{M_2}$  with  $M_1 \in \mathbb{R}^{s \times n_1}$  and  $M_2 \in \mathbb{R}^{s \times n_2}$ , if  $rank(\binom{M_1 \ M_2}{M_1 \ M_2}) = n_1 + n_2$ , then  $f_{M_1}$  and  $f_{M_2}$  are injective and  $Res(f_{M_1}) \cap Res(f_{M_2}) = \emptyset$ .

We prove Proposition 3 in Appendix E. Proposition 1, 2 and 3 provide a new perspective for building powerful aggregators and even injective aggregators. Compared with the distinguishing strength studies in (Xu et al., 2019) and (Corso et al., 2020), as well as existing strategies for building injective aggregators, e.g., sum-of-power mapping (Zaheer et al., 2017) and normalized moments (Corso et al., 2020), we reformulate the aggregation with aggregation coefficients matrix and show the relations of the distinguishing strength of aggregators and the rank of the corresponding aggregation coefficients matrices. Besides, the aggregation of this method is controlled by aggregation coefficients which can be learned from graph data to better leverage structural information. In this paper, to simplify the analysis, we only consider the aggregations within one-hop neighbors. The results can be easily extended to more sophisticated aggregators with the overall framework unchanged

In the perspective of preserving the rank of hidden features among neighbors,  $r = Mh$  indicates that  $rank(r) \leq \min(rank(M), rank(h))$ . To preserve the rank of hidden features in aggregations such that  $rank(r) = rank(h)$ , we need  $rank(M) \geq rank(h)$ . This builds a connection between improving the distinguishing strength of aggregators and preserving the rank of hidden features among neighbors, both of which have the requirements on the rank of  $M$ . General aggregators such as ones in GCN, GIN-0 and GAT have  $rank(M) = 1$ . Thus,  $rank(r)$  is always fixed to 1 no matter what the rank of the input features is. Correspondingly, they have a weak distinguishing strength.

Equation 2 splits the aggregation and feature/structure extraction into two independent steps, which helps to figure out that the expressive power loss happens in the aggregation step, and then the model extracts feature/structure information on the distorted encodings of neighbors. From Equation 2, the aggregation can be considered as a representation regularization step, which unifies different multisets of neighbors into the same representation style while holding permutation invariance. Then, the model can extract structural information on this regulated data with a shared trainable matrix as the third step in Equation 2. Based on this observation, we propose two novel GNN layers: ExpandingConv and CombConv.

### 3.3 EXPANDINGCONV

In this section, we first present ExpandingConv framework. Then we provide one of its implementations and analyze how ExpandingConv achieves more powerful aggregations. The ExpandingConv framework is

$$\begin{cases} \mathbf{m}_{uv}^{(t)} = \mathbf{f}_{\text{local}}(u, v)|_{u \in \mathcal{N}(v)} \\ \mathbf{h}_v^{(t)} = \mathbf{f}_{\text{aggr}}(\{\{\text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T})|u \in \mathcal{N}(v)\}\}), \end{cases}$$

where  $\mathbf{m}_{uv}^{(t)} \in \mathbb{R}^{s \times 1}$  with  $s > 1$  and  $\mathbf{f}_{\text{local}}(u, v)$  is the mapping of local structures between nodes  $u$  and  $v$ . The implementation of  $\mathbf{f}_{\text{local}}(u, v)$  is very flexible with the only restriction of ensuring the same encoding of isomorphic graphs.  $\text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) \in \mathbb{R}^{sd \times 1}$  is the expanded representation of hidden features  $\mathbf{h}_u^{(t-1)} \in \mathbb{R}^{d \times 1}$ . Then a GNN layer  $\mathbf{f}_{\text{aggr}}$  learns structural information on this

expanded representations. We introduce an implementation as follows:

$$\begin{cases} \mathbf{m}_{uv}^{(t)} = \text{Tanh}(\mathbf{W}[\mathbf{h}_v^{(t-1)} \parallel \mathbf{h}_u^{(t-1)}] + \mathbf{b}) \\ \mathbf{h}_v^{(t)} = \sum_{u \in \mathcal{N}(v)} \text{MLP}(\text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T})). \end{cases} \quad (3)$$

In Equation 3, we implement  $f_{\text{local}}(u, v)$  as the function of hidden features of nodes  $u$  and  $v$ . There can be other implementations, and we leave them for future work.  $\mathbf{W} \in \mathbb{R}^{s \times 2d}$  and  $\mathbf{b} \in \mathbb{R}^{s \times 1}$  are trainable matrices. (Luan et al., 2019) empirically showed that different nonlinear activation functions have different contributions in preserving the rank of matrices. We use the recommended Tanh as the activation function in the computation of  $\mathbf{m}_{uv}^{(t)}$  to better preserve the rank of aggregation coefficient matrices. MLP denotes a 2-layer perceptron.

Next, we represent Equation 3 with the corresponding three-stage representation as given in Section 3.2 to obtain its aggregation coefficient matrix and analyze its distinguishing strength. To simplify this process, we only consider 1-layer MLP with  $\mathbf{W}' \in \mathbb{R}^{d \times sd}$  and  $\mathbf{b}' \in \mathbb{R}^{d \times 1}$ .

$$\begin{aligned} \mathbf{h}_v^{(t)} &= \sum_{u \in \mathcal{N}(v)} \text{ReLU}(\mathbf{W}' \text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) + \mathbf{b}') = \\ &= \begin{pmatrix} \sum_{u \in \mathcal{N}(v)} \text{ReLU}(\mathbf{W}'_{[1,:]} \text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) + \mathbf{b}'_{[1]}) \\ \sum_{u \in \mathcal{N}(v)} \text{ReLU}(\mathbf{W}'_{[2,:]} \text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) + \mathbf{b}'_{[2]}) \\ \vdots \\ \sum_{u \in \mathcal{N}(v)} \text{ReLU}(\mathbf{W}'_{[d,:]} \text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) + \mathbf{b}'_{[d]}) \end{pmatrix} = \begin{pmatrix} \sum_{u \in \mathcal{N}_1(v)} (\mathbf{W}'_{[1,:]} \text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) + \mathbf{b}'_{[1]}) \\ \sum_{u \in \mathcal{N}_2(v)} (\mathbf{W}'_{[2,:]} \text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) + \mathbf{b}'_{[2]}) \\ \vdots \\ \sum_{u \in \mathcal{N}_d(v)} (\mathbf{W}'_{[d,:]} \text{vec}(\mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) + \mathbf{b}'_{[d]}) \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{W}'_{[1,:]} & & & \\ & \mathbf{W}'_{[2,:]} & & \\ & & \ddots & \\ & & & \mathbf{W}'_{[d,:]} \end{pmatrix} \begin{pmatrix} \text{vec}(\sum_{u \in \mathcal{N}_1(v)} \mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) \\ \text{vec}(\sum_{u \in \mathcal{N}_2(v)} \mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) \\ \vdots \\ \text{vec}(\sum_{u \in \mathcal{N}_d(v)} \mathbf{m}_{uv}^{(t)} \mathbf{h}_u^{(t-1)T}) \end{pmatrix} + \begin{pmatrix} |\mathcal{N}_1(v)| \cdot \mathbf{b}'_{[1]} \\ |\mathcal{N}_2(v)| \cdot \mathbf{b}'_{[2]} \\ \vdots \\ |\mathcal{N}_d(v)| \cdot \mathbf{b}'_{[d]} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{W}'_{[1,:]} & & & \\ & \mathbf{W}'_{[2,:]} & & \\ & & \ddots & \\ & & & \mathbf{W}'_{[d,:]} \end{pmatrix} \begin{pmatrix} \text{vec}(\mathbf{M}_{v_1 \pi_1}^{(t)} \bar{\mathbf{h}}_{v_1 \pi_1}^{(t-1)T}) \\ \text{vec}(\mathbf{M}_{v_2 \pi_2}^{(t)} \bar{\mathbf{h}}_{v_2 \pi_2}^{(t-1)T}) \\ \vdots \\ \text{vec}(\mathbf{M}_{v_d \pi_d}^{(t)} \bar{\mathbf{h}}_{v_d \pi_d}^{(t-1)T}) \end{pmatrix} + \begin{pmatrix} |\mathcal{N}_1(v)| \cdot \mathbf{b}'_{[1]} \\ |\mathcal{N}_2(v)| \cdot \mathbf{b}'_{[2]} \\ \vdots \\ |\mathcal{N}_d(v)| \cdot \mathbf{b}'_{[d]} \end{pmatrix}, \end{aligned} \quad (4)$$

where  $\mathcal{N}_i(v) \subseteq \mathcal{N}(v) | i \in [d]$  are sampled subsets of neighbors in each dimension.  $\mathbf{M}_{v_i \pi_i}^{(t)} = (\mathbf{m}_{u_1 v_i}^{(t)}, \mathbf{m}_{u_2 v_i}^{(t)}, \dots, \mathbf{m}_{u_{|\mathcal{N}_i(v)|} v_i}^{(t)}) \in \mathbb{R}^{s \times |\mathcal{N}_i(v)|}$  and  $\bar{\mathbf{h}}_{v_i}^{(t-1)} = (\mathbf{h}_{u_1}^{(t-1)}, \mathbf{h}_{u_2}^{(t-1)}, \dots, \mathbf{h}_{u_{|\mathcal{N}_i(v)|}}^{(t-1)}) \in \mathbb{R}^{d \times |\mathcal{N}_i(v)|}$  are aggregation coefficients matrix and hidden feature matrix corresponding to the subset of neighbors  $\mathcal{N}_i(v) \subseteq \mathcal{N}(v)$  according to  $\pi_i$ . We denote  $[\mathbf{h}_v^{(t-1)} \parallel \bar{\mathbf{h}}_{v_i}^{(t-1)}] = (\mathbf{h}_v^{(t-1)} \parallel \mathbf{h}_{u_1}^{(t-1)}, \mathbf{h}_v^{(t-1)} \parallel \mathbf{h}_{u_2}^{(t-1)}, \dots, \mathbf{h}_v^{(t-1)} \parallel \mathbf{h}_{u_{|\mathcal{N}_i(v)|}}^{(t-1)}) \in \mathbb{R}^{2d \times |\mathcal{N}_i(v)|}$ , then  $\mathbf{M}_{v_i}^{(t)} = \text{Tanh}(\mathbf{W}[\mathbf{h}_v^{(t-1)} \parallel \bar{\mathbf{h}}_{v_i}^{(t-1)}] + \mathbf{b}) \in \mathbb{R}^{s \times |\mathcal{N}_i(v)|}$ . According to Equation 4, we finally obtain the three-stage representation equivalent to Equation 3.

$$\begin{cases} \mathbf{M}_{v_i}^{(t)} = \text{Tanh}(\mathbf{W}[\mathbf{h}_v^{(t-1)} \parallel \bar{\mathbf{h}}_{v_i}^{(t-1)}] + \mathbf{b}) \\ \bar{\mathbf{r}}_{v_i}^{(t)} = \mathbf{M}_{v_i \pi_i}^{(t)} \bar{\mathbf{h}}_{v_i \pi_i}^{(t-1)T}, \quad i \in [d] \\ \mathbf{h}_v^{(t)} = \text{diag}(\mathbf{W}'_{[1,:]}, \mathbf{W}'_{[2,:]}, \dots, \mathbf{W}'_{[d,:]})(\text{vec}(\bar{\mathbf{r}}_{v_1}^{(t)}), \text{vec}(\bar{\mathbf{r}}_{v_2}^{(t)}), \dots, \text{vec}(\bar{\mathbf{r}}_{v_d}^{(t)}))^T \\ \quad + (|\mathcal{N}_1(v)| \cdot \mathbf{b}'_{[1]}, |\mathcal{N}_2(v)| \cdot \mathbf{b}'_{[2]}, \dots, |\mathcal{N}_d(v)| \cdot \mathbf{b}'_{[d]})^T. \end{cases} \quad (5)$$

According to the computation of  $\bar{\mathbf{r}}_{v_i}^{(t)}$ ,  $\text{rank}(\bar{\mathbf{r}}_{v_i}^{(t)}) \leq \min(s, d, |\mathcal{N}_i(v)|)$ . By configuring a larger  $s$ , we have  $\text{rank}(\bar{\mathbf{r}}_{v_i}^{(t)}) > 1$  with a high probability, which is different from general GNNs with a rank of 1. As analyzed in Section 3.2, this achieves more powerful aggregators as well as preserves the rank of hidden features among neighbors. The obtained  $\bar{\mathbf{r}}_{v_i}^{(t)}$  after aggregation is the unified representations of neighbors. We then use the trainable matrix  $\mathbf{W}' \in \mathbb{R}^{sd \times d}$  to extract feature/structure information. Unlike the aggregation step, the dimensions reduction here (from  $sd$  to  $d$ ) would not cause information loss. This can be explained by the fact that only task-relevant structural information needs to be preserved and passed to the next layer, and it can be embedded in lower dimensions.

### Comparisons with multi-head GAT.

**Proposition 4.** *Multi-head GAT is an implementation of ExpandingConv as follows:*

$$\begin{cases} \alpha_{vu} = \text{softmax}(\text{LeakyReLU}([\text{diag}(\tilde{\mathbf{a}}^{1T}, \tilde{\mathbf{a}}^{2T}, \dots, \tilde{\mathbf{a}}^{KT}) \\ \|\text{diag}(\tilde{\mathbf{a}}'^{1T}, \tilde{\mathbf{a}}'^{2T}, \dots, \tilde{\mathbf{a}}'^{KT})\|[\mathbf{W}\mathbf{h}_v^{(t-1)}\|\mathbf{W}\mathbf{h}_u^{(t-1)}]])) \\ \mathbf{h}_v^{(t)} = \sigma\left(\frac{1}{K}\mathbf{W}\sum_{j\in\mathcal{N}_v}\text{vec}(\alpha_{vu}\mathbf{h}_u^{(t-1)T})\right), \end{cases} \quad (6)$$

where  $\mathbf{W} = \|\|_{k=1}^K \mathbf{W}^k \in \mathbb{R}^{kd \times d}$  is the concatenation of the trainable matrix in all  $K$  heads.

We prove Proposition 4 in Appendix F. Although multi-head GAT is based on attention mechanism, ExpandingConv provides a new perspective to explain its effectiveness. Applying multi-head attention mechanism helps to preserve the rank of hidden features as well as achieve more powerful aggregators. However, the usage of LeakyReLU may be harmful to preserving the rank of the aggregation coefficient matrix (Luan et al., 2019).

GAT as well as most other GNNs (such as GCN, GIN, etc) follows the same pattern that applies nonlinear units after aggregation. According to the analysis in Section 3.1, Equation 3 applies MLP on  $\text{vec}(\mathbf{m}_{uv}^{(t)}\mathbf{h}_u^{(t-1)T})$  before SUM to break the distinguishing strength limitation of SUM. It also produces other interesting results. By reformulating Equation 3 with its three-stage representation as Equation 5, each dimension of hidden features aggregates on a subset of neighbors independently, which corresponds to a kind of dimension-wise neighbor sampling mechanism. We call the modification of applying ReLU ahead of SUM aggregator as *Re-SUM* mechanism. (Mishra et al., 2020) and (Rong et al., 2019) studied dropedge and node masking mechanism on node-level predictions, both of which can be considered as neighbor sampling strategies that have shown their effectiveness in improving the generalization ability of aggregation-based GNNs and are also used as unbiased data augmentation technique for training. Compared with dropedge and node masking, *Re-SUM* realizes a dimension-wise neighbor sampling, and it does not need to manually set the sampling ratio since this mechanism takes effects implicitly. *Re-SUM* shows that the neural network itself can perform sampling by properly combining nonlinear units and aggregators, without explicitly modifying the network architecture. Our experimental results verified the effectiveness of the *Re-SUM* on a variety of graph tasks.

### 3.4 COMB CONV

The CombConv framework is

$$\begin{cases} \mathbf{m}_{uv}^{(t)} = \mathbf{f}_{\text{local}}(u, v)|_{u \in \mathcal{N}(v)} \\ \mathbf{h}_v^{(t)} = \mathbf{f}_{\text{aggr}}(\{\{\text{vec}(\mathbf{m}_{uv}^{(t)} \odot \mathbf{h}_u^{(t-1)})|_{u \in \mathcal{N}(v)}\}\}), \end{cases}$$

where  $\mathbf{m}_{uv}^{(t)} \in \mathbb{R}^{d \times 1}$  and  $\odot$  denotes element-wise product. An implementation of CombConv is given as follows:

$$\begin{cases} \mathbf{m}_{uv}^{(t)} = \text{Tanh}(\mathbf{W}[\mathbf{h}_v^{(t-1)}\|\mathbf{h}_u^{(t-1)}] + \mathbf{b}) \\ \mathbf{h}_v^{(t)} = \sum_{u \in \mathcal{N}(v)} \text{MLP}(\mathbf{m}_{uv}^{(t)} \odot \mathbf{h}_u^{(t-1)}), \end{cases} \quad (7)$$

where  $\mathbf{W} \in \mathbb{R}^{d \times 2d}$  and  $\mathbf{b} \in \mathbb{R}^{d \times 1}$ . Similar to ExpandingConv, CombConv also applies *Re-SUM* aggregation. The difference is that each dimension of hidden features is aggregated with an independent weighted aggregator. ExpandingConv with  $s = 1$  corresponds to a special case of CombConv where all dimensions share the same aggregator. Therefore, the distinguishing strength of CombConv is stronger than ExpandingConv with  $s = 1$ . Meanwhile, CombConv does not expand the hidden features of nodes in aggregation. Hence, it requires fewer parameters.

## 4 EXPERIMENTS

In this section, we evaluate ExpandingConv and CombConv on graph-level prediction tasks on OGB (Weihua Hu, 2020), TU (Kersting et al., 2016; Yanardag & Vishwanathan, 2015) and QM9 (Ramakrishnan et al., 2014; Wu et al., 2018; Ruddigkeit et al., 2012). The code is available at <https://github.com/qslim/epcb-gnns>.

**Configurations.** We use the default dataset splits for OGB. The QM9 dataset is randomly split into 80% train, 10% validation and 10% test as given in (Morris et al., 2019; Maron et al., 2019). For TU dataset, we follow the standard 10-fold cross validation protocol and splits from (Zhang et al., 2018) and report our results following the protocol described in (Xu et al., 2019; Ying et al., 2018). We use the concatenation of hidden features from all layers to compute the entire graph representations (Xu et al., 2018). In our tests, all models are equipped with batch normalization (Ioffe & Szegedy, 2015) on each hidden layer when evaluating on OGB and TU, and are not when evaluating on QM9. All datasets’ descriptions and detailed hyperparameter settings are given in Appendix H.

We first conduct comprehensive ablation studies to evaluate the effectiveness of powerful aggregators and *Re*-SUM mechanism on OGB and QM9 as given in Table 1 and Table 2. Then, we compare the performance of ExpandingConv and CombConv with competitive baselines on all three datasets as given in Table 3 and Table 4 to show their improvements. ExpC-*s* denotes ExpandingConv with  $W \in \mathbb{R}^{s \times s}$ . We use ExpC\* and CombC\* to denote the ExpandingConv and CombConv without *Re*-SUM.

Table 1: Ablation studies on OGB and QM9. Higher is better.

	ogbg-ppa	ogbg-molhiv	ogbg-molpcba	ogbg-code
ExpC*-1	70.65	77.63	22.65	32.2
ExpC-1	77.50	76.79	23.39	32.6
ExpC-3,4,5	<b>80.11</b>	<b>77.89</b>	23.44	<b>33.2</b>
CombC*	73.61	76.47	23.45	32.29
CombC	77.64	76.63	<b>23.73</b>	32.72

Table 2: Ablation studies on QM9. Lower is better.

	$\mu$	$\alpha$	$\epsilon_{homo}$	$\epsilon_{lumo}$	$\Delta\epsilon$	$\langle R^2 \rangle$	ZPVE	$U_0$	$U$	$H$	$G$	$C_v$
ExpC*-1	0.467	0.283	0.00337	0.00340	0.00467	22.9	0.000205	0.0255	0.0263	0.0242	0.0261	0.1189
ExpC-1	0.469	0.268	0.00326	0.00329	0.00466	20.8	0.000186	0.0202	0.0199	0.0202	0.0201	0.1039
ExpC-4	0.413	0.255	0.00273	0.00300	0.00420	19.4	<b>0.000168</b>	0.0184	0.0183	0.0178	0.0182	0.1115
ExpC-8	0.400	0.257	0.00259	0.00286	0.00395	18.1	0.000172	0.0158	0.0170	0.0177	0.0184	0.1060
ExpC-16	<b>0.382</b>	0.255	<b>0.00248</b>	<b>0.00268</b>	<b>0.00373</b>	17.2	0.000170	0.0170	0.0174	0.0193	<b>0.0165</b>	0.1043
ExpC-32	<b>0.368</b>	<b>0.244</b>	<b>0.00248</b>	<b>0.00257</b>	<b>0.00364</b>	<b>16.3</b>	0.000174	<b>0.0151</b>	<b>0.0167</b>	<b>0.0165</b>	0.0198	<b>0.0962</b>
CombC*	0.4062	0.248	0.00259	0.00273	0.00387	17.1	0.000170	0.0185	0.0181	0.0164	0.0174	0.1022
CombC	0.399	<b>0.241</b>	0.00261	0.00278	0.00386	<b>15.9</b>	<b>0.000160</b>	<b>0.0144</b>	<b>0.0145</b>	<b>0.0147</b>	<b>0.0140</b>	<b>0.0858</b>

#### 4.1 ABLATION STUDIES

**Effect of powerful aggregators.** For complex graph structures with dense connections or with abundant node/edge features, they would benefit from a higher expressive model to maximumly distinguish different structures and extract relevant structural patterns as the model goes deeper to leverage large receptive fields. This is validated on both QM9 and OGB. We configure  $s = 1, 4, 8, 16, 32$  of ExpC-*s* for all 12 targets of QM9. As we apply a larger  $s$ , the model continuously achieves better performance on most targets. We randomly select  $s = 1, 4$  for ogbg-ppa, ogbg-molhiv and  $s = 1, 5$  for ogbg-code. The results show that applying larger  $s$  gains performance improvements, especially on ogbg-ppa which involves large graphs with dense connections.

**Effect of *Re*-SUM mechanism.** In Table 1 and Table 2, the performance differences between ExpC\*-1 (CombC\*) and ExpC-1 (CombC) show the effectiveness of *Re*-SUM. In our tests, the *Re*-SUM can be extremely powerful on graphs with dense connections such as ogbg-ppa, which is validated on both ExpandingConv (with 6.85% improvements) and CombConv (with 4% improvements). On most targets of QM9, this mechanism also gains improvements. For small graphs with sparse connections such as ogbg-hiv and ogbg-molpcba, the improvements are not very significant.

#### 4.2 COMPARISONS WITH BASELINES

Table 3 and Table 4 show the performance comparisons of our models with baselines on QM9, TU and OGB respectively. All datasets in QM9 and OGB graph-level predictions are used for evaluations. For TU, we use 3 widely used datasets: COLLAB includes graphs with dense connections;



Table 3: Comparisons with baselines on OGB and TU. Higher is better.

	OGB				TU		
	ogbg-ppa	ogbg-molhiv	ogbg-molpcba	ogbg-code	COLLAB	RDT-B	RDT-M12
DGK (Yanardag & Vishwanathan, 2015)	NA	NA	NA	NA	73.09 ± 0.25	78.04 ± 0.39	32.22 ± 0.1
PSCN (Niepert et al., 2016)	NA	NA	NA	NA	73.76 ± 0.50	86.30 ± 1.58	41.32 ± 0.42
AWE (Ivanov & Burnaev, 2018)	NA	NA	NA	NA	73.93 ± 1.94	87.89 ± 2.53	39.20 ± 2.09
GCN (Kipf & Welling, 2016)	68.39 ± 0.84	76.06 ± 0.97	20.20 ± 0.24	31.63 ± 0.18	NA	NA	NA
GIN (Xu et al., 2019)	68.92 ± 1.0	75.58 ± 1.40	22.66 ± 0.28	31.63 ± 0.20	80.2 ± 1.9	92.4 ± 2.5	NA
GraphSAG(Hamilton et al., 2017)	NA	NA	NA	NA	68.25	NA	42.24
DiffPool (Ying et al., 2018)	NA	NA	NA	NA	75.48	NA	47.08
CapsGNN (Xinyi & Chen, 2019)	NA	NA	NA	NA	79.62 ± 0.91	NA	46.62 ± 1.9
PPGN (Maron et al., 2019)	NA	NA	NA	NA	80.16 ± 1.1	NA	NA
DeeperGCN (Li et al., 2020a)	77.12 ± 0.71	78.58 ± 1.17	NA	NA	NA	NA	NA
HIMP (Fey et al., 2020)	NA	<b>78.80 ± 0.82</b>	NA	NA	NA	NA	NA
WEGL (Kolouri et al., 2020)	NA	77.57 ± 1.11	NA	NA	NA	NA	NA
multi-head GAT(Veličković et al., 2017)	NA	75.81	20.10	31.10	NA	NA	NA
ExpC-s	<b>79.76 ± 0.72</b>	77.99 ± 0.82	23.42 ± 0.29	<b>33.18 ± 0.17</b>	<b>82.10 ± 1.60</b>	92.2 ± 1.87	<b>49.91 ± 1.75</b>
CombC	77.81 ± 0.76	77.15 ± 1.32	<b>23.63 ± 0.23</b>	32.76 ± 0.15	81.90 ± 1.75	<b>92.5 ± 1.69</b>	49.02 ± 1.21

Table 4: Comparisons with baselines on QM9. Lower is better.

	$\mu$	$\alpha$	$\epsilon_{homo}$	$\epsilon_{lumo}$	$\Delta\epsilon$	$\langle R^2 \rangle$	ZPVE	$U_0$	$U$	$H$	$G$	$C_v$
DTNN (Wu et al., 2018)	<b>0.244</b>	0.95	0.00388	0.00512	0.0112	17	0.00172	2.43	2.43	2.43	2.43	0.27
MPNN (Gilmer et al., 2017)	0.358	0.89	0.00541	0.00623	0.0066	28.5	0.00216	2.05	2	2.02	2.02	0.42
k-GNN (Morris et al., 2019)	0.476	0.27	0.00337	0.00351	0.0048	22.9	0.00019	0.0427	0.111	0.0419	0.0469	<b>0.0944</b>
PPGN (Maron et al., 2019)	<b>0.0934</b>	0.318	<b>0.00174</b>	<b>0.0021</b>	<b>0.0029</b>	<b>3.78</b>	0.000399	0.022	0.0504	0.0294	0.024	0.144
GIN0* (Xu et al., 2019)	0.471	0.281	0.00327	0.00340	0.00473	22.9	0.000202	0.0244	0.0245	0.0233	0.0255	0.1283
GAT-s(Veličković et al., 2017)	0.452	0.286	0.00322	0.00327	0.00460	22.7	0.000228	0.0212	0.0223	0.0223	0.0219	0.1247
ExpC-s	0.368	<b>0.244</b>	<b>0.00248</b>	<b>0.00257</b>	<b>0.00364</b>	16.3	<b>0.000168</b>	<b>0.0151</b>	<b>0.0167</b>	<b>0.0165</b>	<b>0.0165</b>	0.0962
CombC	0.399	<b>0.241</b>	0.00261	0.00278	0.00386	<b>15.9</b>	<b>0.000160</b>	<b>0.0144</b>	<b>0.0145</b>	<b>0.0147</b>	<b>0.0140</b>	<b>0.0858</b>

REDDIT-BINARY (RDT-B) and REDDIT-MULTI-12K (RDT-M12) include large and sparse graphs with one center node having dense connections with other nodes. All results of baselines are taken from the original papers except for the results of GraphSAGE on TU, multi-head GAT on OGB and GIN0\* on QM9 which were not reported by the original papers. We report the results of GraphSAGE provided by (Ying et al., 2018) and evaluate multi-head GAT and GIN0\* by ourselves. To ensure a fair comparison, for OGB and TU, we configure the number of heads in multi-head GAT and  $s$  in ExpC-s to be the same which is selected in  $\{3, 4, 5\}$ . For QM9, the number of heads is 8 and  $s \in \{4, 8, 16, 32\}$ . GIN0\* in QM9 denotes GIN0 without batch normalization.

Compared with baselines, our models achieve the best performance on 7 out of all 12 targets of QM9, 3 out of all 4 graph-level prediction datasets of OGB and all 3 selected TU datasets. Our models get 1.9% improvements on COLLAB and 2.83% improvements on REDDIT-MULTI-12K compared with SOTA baselines. On ogbg-ppa, our models achieve 2.6% higher classification accuracies compared with SOTA baselines. On ogbg-code, they achieve 1.5% improvements. Multi-head GAT can also be considered as an implementation of ExpandingConv. However, its performance on graph-level predictions is not competitive. According to its three-stage representation, the usage of LeakyReLU in the aggregation step is harmful to preserving the rank, and the usage of softmax makes it harder to analyze the rank. In the extraction step, the 1-layer MLP may have a limited representation power to represent the desired extraction functions.

## 5 CONCLUSION

We show how basic aggregators used in general GNNs become expressive bottlenecks. To address this limitation, we develop theoretical foundations of building powerful aggregators. We also propose the *Re-SUM* mechanism which achieves dimension-wise sampling. To evaluate their effectiveness, we develop two novel GNN layers, and conduct extensive experiments on public graph benchmarks. The results are consistent with our analysis, and our proposed models achieve SOTA performance on a variety of graph-level prediction benchmarks.

## REFERENCES

Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.

- Zhengdao Chen, Soledad Villar, Lei Chen, and Joan Bruna. On the equivalence between graph isomorphism testing and function approximation with gnns. In *Advances in Neural Information Processing Systems*, pp. 15868–15876, 2019.
- Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. *arXiv preprint arXiv:2004.05718*, 2020.
- Nima Dehmamy, Albert-László Barabási, and Rose Yu. Understanding the representation power of graph neural networks in learning graph topology. In *Advances in Neural Information Processing Systems*, pp. 15413–15423, 2019.
- David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in neural information processing systems*, pp. 2224–2232, 2015.
- Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Benchmarking graph neural networks. *arXiv preprint arXiv:2003.00982*, 2020.
- M. Fey, J. G. Yuen, and F. Weichert. Hierarchical inter-message passing for learning on molecular graphs. In *ICML Graph Representation Learning and Beyond (GRL+) Workshop*, 2020.
- Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pp. 1263–1272. JMLR. org, 2017.
- Marco Gori, Gabriele Monfardini, and Franco Scarselli. A new model for learning in graph domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 2, pp. 729–734. IEEE, 2005.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pp. 1024–1034, 2017.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Sergey Ivanov and Evgeny Burnaev. Anonymous walk embeddings. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2191–2200, Stockholm, Sweden, 10–15 Jul 2018. PMLR. URL <http://proceedings.mlr.press/v80/ivanov18a.html>.
- Kristian Kersting, Nils M. Kriege, Christopher Morris, Petra Mutzel, and Marion Neumann. Benchmark data sets for graph kernels, 2016. <http://graphkernels.cs.tu-dortmund.de>.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Soheil Kolouri, Navid Naderializadeh, Gustavo K Rohde, and Heiko Hoffmann. Wasserstein embedding for graph learning. *arXiv preprint arXiv:2006.09430*, 2020.
- Guohao Li, Chenxin Xiong, Ali Thabet, and Bernard Ghanem. Deepergcn: All you need to train deeper gnns. *arXiv preprint arXiv:2006.07739*, 2020a.
- Pan Li, Yanbang Wang, Hongwei Wang, and Jure Leskovec. Distance encoding–design provably more powerful gnns for structural representation learning. *arXiv preprint arXiv:2009.00142*, 2020b.
- Yaw-Ling Lin and Steven S Skiena. Algorithms for square roots of graphs. *SIAM Journal on Discrete Mathematics*, 8(1):99–118, 1995.

- Sitao Luan, Mingde Zhao, Xiao-Wen Chang, and Doina Precup. Break the ceiling: Stronger multi-scale deep graph convolutional networks. In *Advances in neural information processing systems*, pp. 10945–10955, 2019.
- Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph networks. *arXiv preprint arXiv:1905.11136*, 2019.
- Pushkar Mishra, Aleksandra Piktus, Gerard Goossen, and Fabrizio Silvestri. Node masking: Making graph neural networks generalize and scale better. *arXiv preprint arXiv:2001.07524*, 2020.
- Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 4602–4609, 2019.
- Ryan L Murphy, Balasubramaniam Srinivasan, Vinayak Rao, and Bruno Ribeiro. Janossy pooling: Learning deep permutation-invariant functions for variable-size inputs. *arXiv preprint arXiv:1811.01900*, 2018.
- Ryan L Murphy, Balasubramaniam Srinivasan, Vinayak Rao, and Bruno Ribeiro. Relational pooling for graph representations. *arXiv preprint arXiv:1903.02541*, 2019.
- Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. Learning convolutional neural networks for graphs. In *International Conference on Machine Learning*, pp. 2014–2023, 2016.
- Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1:140022, 2014.
- Yu Rong, Wenbing Huang, Tingyang Xu, and Junzhou Huang. Dropedge: Towards deep graph convolutional networks on node classification. In *International Conference on Learning Representations*, 2019.
- Lars Ruddigkeit, Ruud Van Deursen, Lorenz C Blum, and Jean-Louis Reymond. Enumeration of 166 billion organic small molecules in the chemical universe database gdb-17. *Journal of chemical information and modeling*, 52(11):2864–2875, 2012.
- Ryoma Sato. A survey on the expressive power of graph neural networks. *arXiv preprint arXiv:2003.04078*, 2020.
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2008.
- Gunther Schmidt. *Relational mathematics*, volume 132. Cambridge University Press, 2011.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- Clement Vignac, Andreas Loukas, and Pascal Frossard. Building powerful and equivariant graph neural networks with message-passing. *arXiv preprint arXiv:2006.15107*, 2020.
- Marinka Zitnik Yuxiao Dong Hongyu Ren Bowen Liu Michele Catasta Jure Leskovec Weihua Hu, Matthias Fey. Open graph benchmark: Datasets for machine learning on graphs. *arXiv preprint arXiv:2005.00687*, 2020.
- B. Yu. Weisfeiler and A. A. Leman. Reduction of a graph to a canonical form and an algebra arising during this reduction. 1968.
- Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):513–530, 2018.
- Zhang Xinyi and Lihui Chen. Capsule graph neural network. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Byl8BnRcYm>.

- Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. *arXiv preprint arXiv:1806.03536*, 2018.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryGs6iA5Km>.
- Pinar Yanardag and SVN Vishwanathan. Deep graph kernels. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1365–1374. ACM, 2015.
- Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. In *Advances in Neural Information Processing Systems*, pp. 4800–4810, 2018.
- Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In *Advances in neural information processing systems*, pp. 3391–3401, 2017.
- Muhan Zhang, Zhicheng Cui, Marion Neumann, and Yixin Chen. An end-to-end deep learning architecture for graph classification. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

## A GCN, GAT AND GIN

Here, we present implementations of GCN, GAT and GIN for the usage of our analysis.

**Graph Convolution Networks (GCN) (Kipf & Welling, 2016).**

$$\mathbf{H}^{(t)} = \sigma(\hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(t-1)} \mathbf{W} + \mathbf{b}) \quad (8)$$

**Graph Attention Networks (GAT) (Veličković et al., 2017).**

$$\mathbf{h}_v^{(t)} = \sigma \left( \sum_{u \in \mathcal{N}(v)} \text{att}(\mathbf{h}_v^{(t-1)}, \mathbf{h}_u^{(t-1)}) \mathbf{W} \mathbf{h}_u^{(t-1)} \right) \quad (9)$$

**Graph Isomorphism Networks (GIN-0) (Xu et al., 2019).**

$$\mathbf{h}_v^{(t)} = \text{MLP} \left( \sum_{u \in \mathcal{N}(v)} \mathbf{h}_u^{(t-1)} \right). \quad (10)$$

## B PROOF OF LEMMA 1

*Proof.* (i) For any two multisets  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , if  $g(\mathbf{f}_{\text{aggr}}(\mathbf{x}_1)) \neq g(\mathbf{f}_{\text{aggr}}(\mathbf{x}_2))$ , then  $\mathbf{f}_{\text{aggr}}(\mathbf{x}_1) \neq \mathbf{f}_{\text{aggr}}(\mathbf{x}_2)$ . Therefore, we have  $\mathbf{f}_{\text{aggr}} \geq g \circ \mathbf{f}_{\text{aggr}}$ . If  $g$  is injective, then  $\mathbf{f}_{\text{aggr}}(\mathbf{x}_1) \neq \mathbf{f}_{\text{aggr}}(\mathbf{x}_2) \Rightarrow g(\mathbf{f}_{\text{aggr}}(\mathbf{x}_1)) \neq g(\mathbf{f}_{\text{aggr}}(\mathbf{x}_2))$ . We have  $\mathbf{f}_{\text{aggr}} \geq g \circ \mathbf{f}_{\text{aggr}} \geq \mathbf{f}_{\text{aggr}}$ , therefore  $\mathbf{f}_{\text{aggr}} = g \circ \mathbf{f}_{\text{aggr}}$ .

(ii) For any two multisets  $\mathbf{x}_1$  and  $\mathbf{x}_2$ ,  $\mathbf{f}_{\text{aggr}1}(\mathbf{x}_1) \neq \mathbf{f}_{\text{aggr}1}(\mathbf{x}_2) \Rightarrow [\mathbf{f}_{\text{aggr}1}(\mathbf{x}_1) \parallel \mathbf{f}_{\text{aggr}2}(\mathbf{x}_1)] \neq [\mathbf{f}_{\text{aggr}1}(\mathbf{x}_2) \parallel \mathbf{f}_{\text{aggr}2}(\mathbf{x}_2)]$ . Therefore,  $\mathbf{f}_{\text{aggr}1} \otimes \mathbf{f}_{\text{aggr}2} \geq \mathbf{f}_{\text{aggr}1}$ . If  $\mathbf{f}_{\text{aggr}1}$  and  $\mathbf{f}_{\text{aggr}2}$  are incomparable, there exist  $\mathbf{x}_3$  and  $\mathbf{x}_4$  such that  $\mathbf{f}_{\text{aggr}2}(\mathbf{x}_3) \neq \mathbf{f}_{\text{aggr}2}(\mathbf{x}_4)$  but  $\mathbf{f}_{\text{aggr}1}(\mathbf{x}_3) = \mathbf{f}_{\text{aggr}1}(\mathbf{x}_4)$ . Therefore, there exist  $\mathbf{x}_3$  and  $\mathbf{x}_4$  such that  $[\mathbf{f}_{\text{aggr}1}(\mathbf{x}_3) \parallel \mathbf{f}_{\text{aggr}2}(\mathbf{x}_3)] \neq [\mathbf{f}_{\text{aggr}1}(\mathbf{x}_4) \parallel \mathbf{f}_{\text{aggr}2}(\mathbf{x}_4)]$  but  $\mathbf{f}_{\text{aggr}1}(\mathbf{x}_3) = \mathbf{f}_{\text{aggr}1}(\mathbf{x}_4)$ .  $\mathbf{f}_{\text{aggr}1} \otimes \mathbf{f}_{\text{aggr}2} \succ \mathbf{f}_{\text{aggr}1}$ .

(iii) Since  $\mathbf{f}_{\text{aggr}}$  is an equivariant aggregator, then  $\mathbf{f}_{\text{aggr}}(\mathbf{T} \cdot \mathbf{x}_1, \mathbf{T} \cdot \mathbf{x}_2, \dots, \mathbf{T} \cdot \mathbf{x}_n) = \mathbf{T} \cdot \mathbf{f}_{\text{aggr}}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \leq \mathbf{f}_{\text{aggr}}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ .  $\square$

## C PROOF OF PROPOSITION 1

*Proof.* (i)

$$\mathbf{f}_{\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)}(\mathbf{x}) = \left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)_\pi \mathbf{x}_\pi = \begin{pmatrix} \mathbf{m}_\pi \mathbf{x}_\pi \\ \mathbf{m}'_\pi \mathbf{x}_\pi \end{pmatrix} = \begin{pmatrix} \mathbf{f}_M(\mathbf{x}) \\ \mathbf{f}_{M'}(\mathbf{x}) \end{pmatrix},$$

then for any  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , we have

$$\mathbf{f}_M(\mathbf{x}_1) \neq \mathbf{f}_M(\mathbf{x}_2) \Rightarrow \mathbf{f}_{\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)}(\mathbf{x}_1) \neq \mathbf{f}_{\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)}(\mathbf{x}_2),$$

and therefore we conclude that  $\mathbf{f}_{\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)} \geq \mathbf{f}_M$ .

(ii) “ $\mathbf{f}_{\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)} > \mathbf{f}_M \Leftrightarrow \text{rank}\left(\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)\right) > \text{rank}(M)$ ”

We prove the claim by contradiction. Assume that  $\text{rank}\left(\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)\right) > \text{rank}(M)$  and  $\mathbf{f}_{\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)} = \mathbf{f}_M$ .  $\mathbf{f}_{\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)} = \mathbf{f}_M$  means that for any  $\mathbf{x}_1$  and  $\mathbf{x}_2$ ,  $\mathbf{m}_\pi \mathbf{x}_{1\pi} = \mathbf{m}_\pi \mathbf{x}_{2\pi} \Leftrightarrow \left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)_\pi \mathbf{x}_{1\pi} = \left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)_\pi \mathbf{x}_{2\pi}$ , where  $\pi$  is the ordering of input elements. Let  $\mathbf{s} = \mathbf{x}_{1\pi} - \mathbf{x}_{2\pi}$ . For any  $i \in [n]$ ,  $s[i] = x_{1\pi}[i] - x_{2\pi}[i] \in \mathbb{R}$ . Then for any  $\mathbf{s} \in \mathbb{R}^n$ ,  $\mathbf{m}_\pi \mathbf{s} = \mathbf{0} \Leftrightarrow \left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)_\pi \mathbf{s} = \mathbf{0}$ . The system of linear equations  $\mathbf{m}_\pi \mathbf{x} = \mathbf{0}$  and  $\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)_\pi \mathbf{x} = \mathbf{0}$  share the same solution space. Let  $R_S$  denote the rank of this solution space, then  $\text{rank}(\mathbf{m}_\pi) + R_S = \text{rank}\left(\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)_\pi\right) + R_S = n$ . Therefore,  $\text{rank}(\mathbf{m}_\pi) = \text{rank}\left(\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)_\pi\right)$ , then we have  $\text{rank}(M) = \text{rank}\left(\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)\right)$ . Since we assumed that  $\text{rank}\left(\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)\right) > \text{rank}(M)$ , we reach a contradiction.

“ $\mathbf{f}_{\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)} > \mathbf{f}_M \Rightarrow \text{rank}\left(\left(\begin{smallmatrix} M \\ M' \end{smallmatrix}\right)\right) > \text{rank}(M)$ ”

We prove an equivalent proposition “ $\text{rank}(\begin{pmatrix} M \\ M' \end{pmatrix}) \leq \text{rank}(M) \Rightarrow \mathbf{f}_{\begin{pmatrix} M \\ M' \end{pmatrix}} \leq \mathbf{f}_M$ ”. Note that  $\text{rank}(\begin{pmatrix} M \\ M' \end{pmatrix}) \geq \text{rank}(M)$  and  $\mathbf{f}_{\begin{pmatrix} M \\ M' \end{pmatrix}} \geq \mathbf{f}_M$  as given in Proposition 1(i). We only need to prove “ $\text{rank}(\begin{pmatrix} M \\ M' \end{pmatrix}) = \text{rank}(M) \Rightarrow \mathbf{f}_{\begin{pmatrix} M \\ M' \end{pmatrix}} = \mathbf{f}_M$ ”.  $\text{rank}(\begin{pmatrix} M \\ M' \end{pmatrix}) = \text{rank}(M)$  means that any row in  $M'$  is linearly dependent to rows in  $M$ . Therefore, there exists  $L \in \mathbb{R}^{s' \times s}$  so that  $\begin{pmatrix} M \\ M' \end{pmatrix} = \begin{pmatrix} I \\ L \end{pmatrix} M$ . For any  $x_1$  and  $x_2$  with  $M P_\pi x_{1\pi} = M P_\pi x_{2\pi}$ ,  $\begin{pmatrix} I \\ L \end{pmatrix} M P_\pi x_{1\pi} = \begin{pmatrix} I \\ L \end{pmatrix} M P_\pi x_{2\pi}$ , and therefore  $\begin{pmatrix} M \\ M' \end{pmatrix} P_\pi x_{1\pi} = \begin{pmatrix} M \\ M' \end{pmatrix} P_\pi x_{2\pi}$ , where  $\pi$  is the ordering of input elements. That is, for any  $x_1$  and  $x_2$ ,  $\mathbf{f}_M(x_1) = \mathbf{f}_M(x_2) \Rightarrow \mathbf{f}_{\begin{pmatrix} M \\ M' \end{pmatrix}}(x_1) = \mathbf{f}_{\begin{pmatrix} M \\ M' \end{pmatrix}}(x_2)$ , thus  $\mathbf{f}_{\begin{pmatrix} M \\ M' \end{pmatrix}} \leq \mathbf{f}_M$ . Finally, we have  $\mathbf{f}_{\begin{pmatrix} M \\ M' \end{pmatrix}} = \mathbf{f}_M$ .

(iii) “Any multiset of size  $n$  is distinguishable with  $\mathbf{f}_M \Rightarrow \text{rank}(M) = n$ ”

Since  $\text{rank}(M) \leq n$ , we prove an equivalent proposition “ $\text{rank}(M) < n \Rightarrow$  there exists at least two multisets which are indistinguishable”. Considering the system of linear equations  $\mathbf{y} = M\mathbf{x}$  where  $\mathbf{x} \in \mathbb{R}^n$ , if  $\text{rank}(M) < n$ , then there exists  $\mathbf{y}'$  such that  $\text{rank}(M) = \text{rank}(M, \mathbf{y}') < n$ . According to the Rouché–Capelli theorem, there are infinite solutions  $\mathbf{x}'_i$  such that  $\mathbf{y}' = M\mathbf{x}'_1 = M\mathbf{x}'_2 = \dots = M\mathbf{x}'_\infty$ . Each  $\mathbf{x}'_i$  comes from a multiset with a particular order. Next, we need to prove that all these  $\mathbf{x}'_i$  come from more than one multiset. As a multiset with bounded size  $n$  constitutes at most  $n!$  different orders, the infinite number of  $\mathbf{x}'_i$  corresponds to  $\mathbf{y}'$  must come from more than one multisets, making these multisets indistinguishable.

“Any multiset of size  $n$  is distinguishable with  $\mathbf{f}_M \Leftarrow \text{rank}(M) = n$ ”

Since  $\text{rank}(M) = n$  and  $s = n$ , for any  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} = M\mathbf{x} \in \mathbb{R}^n$  is unique. Correspondingly, for any  $P_\pi x_\pi$ ,  $M(P_\pi x_\pi)$  is unique. □

## D PROOF OF PROPOSITION 2

*Proof.* (i) According to the proof of Proposition 1(i),  $\mathbf{f}_{\begin{pmatrix} M \\ M' \end{pmatrix}}(\mathbf{x}) = \begin{pmatrix} \mathbf{f}_M(\mathbf{x}) \\ \mathbf{f}_{M'}(\mathbf{x}) \end{pmatrix}$ . For any  $x_1$  and  $x_2$  with  $\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}(x_1) = \mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}}(x_2)$ ,  $\mathbf{f}_{M_1}(x_1) = \mathbf{f}_{M_2}(x_2)$  holds. Meanwhile,  $\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}(x_1) = \mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}}(x_2) \in \text{Res}(\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}) \cap \text{Res}(\mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}})$ , and  $\mathbf{f}_{M_1}(x_1) = \mathbf{f}_{M_2}(x_2) \in \text{Res}(\mathbf{f}_{M_1}) \cap \text{Res}(\mathbf{f}_{M_2})$ . Therefore, for any  $e \in \text{Res}(\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}) \cap \text{Res}(\mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}})$ , we have  $e \in \text{Res}(\mathbf{f}_{M_1}) \cap \text{Res}(\mathbf{f}_{M_2})$ . That is  $\text{Res}(\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}) \cap \text{Res}(\mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}}) \subseteq \text{Res}(\mathbf{f}_{M_1}) \cap \text{Res}(\mathbf{f}_{M_2})$ .

(ii) We prove an equivalent proposition “ $\text{Res}(\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}) \cap \text{Res}(\mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}}) = \text{Res}(\mathbf{f}_{M_1}) \cap \text{Res}(\mathbf{f}_{M_2}) \Leftarrow \text{rank}(\begin{pmatrix} M_1 & M_2 \\ M'_1 & M'_2 \end{pmatrix}) = \text{rank}(\begin{pmatrix} M_1 & M_2 \end{pmatrix})$ ”.  $\text{rank}(\begin{pmatrix} M_1 & M_2 \\ M'_1 & M'_2 \end{pmatrix}) = \text{rank}(\begin{pmatrix} M_1 & M_2 \end{pmatrix})$  means that any row in  $\begin{pmatrix} M_1 & M_2 \\ M'_1 & M'_2 \end{pmatrix}$  is linearly dependent to  $\begin{pmatrix} M_1 & M_2 \end{pmatrix}$ . Therefore, there exists  $L \in \mathbb{R}^{s' \times s}$  so that  $\begin{pmatrix} M_1 & M_2 \\ M'_1 & M'_2 \end{pmatrix} = \begin{pmatrix} I \\ L \end{pmatrix} \begin{pmatrix} M_1 & M_2 \end{pmatrix}$ . Correspondingly,  $\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix} = \begin{pmatrix} I \\ L \end{pmatrix} M_1$  and  $\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix} = \begin{pmatrix} I \\ L \end{pmatrix} M_2$ . For any  $x_1$  and  $x_2$  with  $M_1 P_\pi x_{1\pi} = M_2 P_\pi x_{2\pi}$ ,  $\begin{pmatrix} I \\ L \end{pmatrix} M_1 P_\pi x_{1\pi} = \begin{pmatrix} I \\ L \end{pmatrix} M_2 P_\pi x_{2\pi}$ , and therefore  $\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix} P_\pi x_{1\pi} = \begin{pmatrix} M_2 \\ M'_2 \end{pmatrix} P_\pi x_{2\pi}$ , where  $\pi$  is the ordering of input elements. Thus for any  $x_1$  and  $x_2$ ,  $\mathbf{f}_{M_1}(x_1) = \mathbf{f}_{M_2}(x_2) \in \text{Res}(\mathbf{f}_{M_1}) \cap \text{Res}(\mathbf{f}_{M_2}) \Rightarrow \mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}(x_1) = \mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}}(x_2) \in \text{Res}(\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}) \cap \text{Res}(\mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}})$ . Hence,  $\text{Res}(\mathbf{f}_{M_1}) \cap \text{Res}(\mathbf{f}_{M_2}) \subseteq \text{Res}(\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}) \cap \text{Res}(\mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}})$ . According to Proposition 2(i),  $\text{Res}(\mathbf{f}_{M_1}) \cap \text{Res}(\mathbf{f}_{M_2}) = \text{Res}(\mathbf{f}_{\begin{pmatrix} M_1 \\ M'_1 \end{pmatrix}}) \cap \text{Res}(\mathbf{f}_{\begin{pmatrix} M_2 \\ M'_2 \end{pmatrix}})$ . □

## E PROOF OF PROPOSITION 3

*Proof.* Since  $M_1 \in \mathbb{R}^{s \times n_1}$ ,  $M_2 \in \mathbb{R}^{s \times n_2}$  and  $\text{rank}(\begin{pmatrix} M_1 & M_2 \end{pmatrix}) = n_1 + n_2$ , we have  $\text{rank}(M_1) = n_1$  and  $\text{rank}(M_2) = n_2$ . According to Proposition 1,  $f_{M_1}$  and  $f_{M_2}$  are injective.

We build the system of linear equations  $\mathbf{y} = \mathbf{A}\mathbf{x}$ , where  $\mathbf{x} \in \mathbb{R}^{n_1+n_2}$ , and  $\mathbf{A} = \begin{pmatrix} M_1 & M_2 \end{pmatrix} \begin{pmatrix} I^{n_1} & \mathbf{0} \\ \mathbf{0} & -I^{n_2} \end{pmatrix} \in \mathbb{R}^{s \times (n_1+n_2)}$ . Then,  $\text{rank}(\mathbf{A}) = \text{rank}(\begin{pmatrix} M_1 & M_2 \end{pmatrix} \begin{pmatrix} I^{n_1} & \mathbf{0} \\ \mathbf{0} & -I^{n_2} \end{pmatrix}) = \text{rank}(\begin{pmatrix} M_1 & M_2 \end{pmatrix}) = n_1 + n_2$ , which means  $\mathbf{A}\mathbf{x} = \mathbf{0}$  has no non-zero solutions. Let  $\mathbf{x}' = (\mathbf{x}[1], \mathbf{x}[2], \dots, \mathbf{x}[n_1])$  and  $\mathbf{x}'' = (\mathbf{x}[n_1 + 1], \mathbf{x}[n_1 + 2], \dots, \mathbf{x}[n_1 + n_2])$  such that  $\mathbf{x} = \begin{pmatrix} \mathbf{x}' \\ \mathbf{x}'' \end{pmatrix}$ . For any  $\mathbf{x} \neq \mathbf{0}$ ,

$$\mathbf{A}\mathbf{x} = \begin{pmatrix} M_1 & M_2 \end{pmatrix} \begin{pmatrix} I^{n_1} & \mathbf{0} \\ \mathbf{0} & -I^{n_2} \end{pmatrix} \begin{pmatrix} \mathbf{x}' \\ \mathbf{x}'' \end{pmatrix} = M_1\mathbf{x}' - M_2\mathbf{x}'' \neq \mathbf{0}.$$

Therefore, for any  $\mathbf{x}' \in \mathbb{R}^{n_1}$ ,  $\mathbf{x}'' \in \mathbb{R}^{n_2}$  and  $\mathbf{x}', \mathbf{x}'' \neq \mathbf{0}$ ,  $M_1\mathbf{x}' \neq M_2\mathbf{x}''$ , hence  $M_1(\mathbf{P}_\pi \mathbf{x}'_\pi) \neq M_2(\mathbf{P}_\pi \mathbf{x}''_\pi)$  for any  $\mathbf{P}_\pi \mathbf{x}'_\pi$  and  $\mathbf{P}_\pi \mathbf{x}''_\pi$ . As a result,  $\text{Res}(f_{M_1}) \cap \text{Res}(f_{M_2}) = \emptyset$ .  $\square$

## F PROOF OF PROPOSITION 4

*Proof.* For Multi-head GAT, there are two types of implementations on aggregating each head, *concatenation* and *average*. Here, we only consider the average aggregation implementation.

$$\begin{aligned} \mathbf{h}_v^{(t)} &= \sigma \left( \frac{1}{K} \sum_{k=1}^K \sum_{j \in \mathcal{N}_v} \alpha_{vu}^k \mathbf{W}^k \mathbf{h}_u^{(t-1)} \right) \\ &= \sigma \left( \frac{1}{K} \sum_{j \in \mathcal{N}_v} \sum_{k=1}^K \mathbf{W}^k (\alpha_{vu}^k \mathbf{h}_u^{(t-1)}) \right) \\ &= \sigma \left( \frac{1}{K} \sum_{j \in \mathcal{N}_v} (\|\mathbf{W}^k\|_{k=1}^K) (\|\alpha_{vu}^k \mathbf{h}_u^{(t-1)}\|_{k=1}^K) \right) \\ &= \sigma \left( \frac{1}{K} \sum_{j \in \mathcal{N}_v} \mathbf{W} \text{vec}(\alpha_{vu} \mathbf{h}_u^{(t-1)T}) \right) \\ &= \sigma \left( \frac{1}{K} \mathbf{W} \sum_{j \in \mathcal{N}_v} \text{vec}(\alpha_{vu} \mathbf{h}_u^{(t-1)T}) \right), \end{aligned}$$

where  $\mathbf{W}^k \in \mathbb{R}^{d \times d}$  is the trainable matrix for the  $k$ -th head, and  $\mathbf{W} = \|\mathbf{W}^k\|_{k=1}^K \in \mathbb{R}^{Kd \times d}$  is the concatenation of the trainable matrix in all  $K$  heads;

$$\alpha_{vu} = \text{softmax}(\text{LeakyReLU} \begin{pmatrix} \mathbf{a}^{1T} [\mathbf{W}^1 \mathbf{h}_v^{(t-1)} \| \mathbf{W}^1 \mathbf{h}_u^{(t-1)}] \\ \mathbf{a}^{2T} [\mathbf{W}^2 \mathbf{h}_v^{(t-1)} \| \mathbf{W}^2 \mathbf{h}_u^{(t-1)}] \\ \vdots \\ \mathbf{a}^{KT} [\mathbf{W}^K \mathbf{h}_v^{(t-1)} \| \mathbf{W}^K \mathbf{h}_u^{(t-1)}] \end{pmatrix}).$$

Let  $\tilde{\mathbf{a}}^* = (\mathbf{a}_1^*, \mathbf{a}_2^*, \dots, \mathbf{a}_K^*)$  and  $\tilde{\mathbf{a}}'^* = (\mathbf{a}_{K+1}^*, \mathbf{a}_{K+2}^*, \dots, \mathbf{a}_{2K}^*)$ . Then,

$$\begin{aligned}
& \begin{pmatrix} \mathbf{a}^{1T} [\mathbf{W}^1 \mathbf{h}_v^{(t-1)} \parallel \mathbf{W}^1 \mathbf{h}_u^{(t-1)}] \\ \mathbf{a}^{2T} [\mathbf{W}^2 \mathbf{h}_v^{(t-1)} \parallel \mathbf{W}^2 \mathbf{h}_u^{(t-1)}] \\ \vdots \\ \mathbf{a}^{KT} [\mathbf{W}^K \mathbf{h}_v^{(t-1)} \parallel \mathbf{W}^K \mathbf{h}_u^{(t-1)}] \end{pmatrix} \\
&= \begin{pmatrix} [\tilde{\mathbf{a}}^1 \parallel \tilde{\mathbf{a}}'^1]^T [\mathbf{W}^1 \mathbf{h}_v^{(t-1)} \parallel \mathbf{W}^1 \mathbf{h}_u^{(t-1)}] \\ [\tilde{\mathbf{a}}^2 \parallel \tilde{\mathbf{a}}'^2]^T [\mathbf{W}^2 \mathbf{h}_v^{(t-1)} \parallel \mathbf{W}^2 \mathbf{h}_u^{(t-1)}] \\ \vdots \\ [\tilde{\mathbf{a}}^K \parallel \tilde{\mathbf{a}}'^K]^T [\mathbf{W}^K \mathbf{h}_v^{(t-1)} \parallel \mathbf{W}^K \mathbf{h}_u^{(t-1)}] \end{pmatrix} \\
&= \begin{pmatrix} \tilde{\mathbf{a}}^{1T} \mathbf{W}^1 \mathbf{h}_v^{(t-1)} + \tilde{\mathbf{a}}'^{1T} \mathbf{W}^1 \mathbf{h}_u^{(t-1)} \\ \tilde{\mathbf{a}}^{2T} \mathbf{W}^2 \mathbf{h}_v^{(t-1)} + \tilde{\mathbf{a}}'^{2T} \mathbf{W}^2 \mathbf{h}_u^{(t-1)} \\ \vdots \\ \tilde{\mathbf{a}}^{KT} \mathbf{W}^K \mathbf{h}_v^{(t-1)} + \tilde{\mathbf{a}}'^{KT} \mathbf{W}^K \mathbf{h}_u^{(t-1)} \end{pmatrix} \\
&= \begin{pmatrix} \tilde{\mathbf{a}}^{1T} \mathbf{W}^1 \mathbf{h}_v^{(t-1)} \\ \tilde{\mathbf{a}}^{2T} \mathbf{W}^2 \mathbf{h}_v^{(t-1)} \\ \vdots \\ \tilde{\mathbf{a}}^{KT} \mathbf{W}^K \mathbf{h}_v^{(t-1)} \end{pmatrix} + \begin{pmatrix} \tilde{\mathbf{a}}'^{1T} \mathbf{W}^1 \mathbf{h}_u^{(t-1)} \\ \tilde{\mathbf{a}}'^{2T} \mathbf{W}^2 \mathbf{h}_u^{(t-1)} \\ \vdots \\ \tilde{\mathbf{a}}'^{KT} \mathbf{W}^K \mathbf{h}_u^{(t-1)} \end{pmatrix} \\
&= \begin{pmatrix} \tilde{\mathbf{a}}^{1T} & & & \\ & \tilde{\mathbf{a}}^{2T} & & \\ & & \ddots & \\ & & & \tilde{\mathbf{a}}^{KT} \end{pmatrix} \begin{pmatrix} \mathbf{W}^1 \\ \mathbf{W}^2 \\ \vdots \\ \mathbf{W}^K \end{pmatrix} \mathbf{h}_v^{(t-1)} + \begin{pmatrix} \tilde{\mathbf{a}}'^{1T} & & & \\ & \tilde{\mathbf{a}}'^{2T} & & \\ & & \ddots & \\ & & & \tilde{\mathbf{a}}'^{KT} \end{pmatrix} \begin{pmatrix} \mathbf{W}^1 \\ \mathbf{W}^2 \\ \vdots \\ \mathbf{W}^K \end{pmatrix} \mathbf{h}_u^{(t-1)} \\
&= \text{diag}(\tilde{\mathbf{a}}^{1T}, \tilde{\mathbf{a}}^{2T}, \dots, \tilde{\mathbf{a}}^{KT}) \mathbf{W} \mathbf{h}_v^{(t-1)} + \text{diag}(\tilde{\mathbf{a}}'^{1T}, \tilde{\mathbf{a}}'^{2T}, \dots, \tilde{\mathbf{a}}'^{KT}) \mathbf{W} \mathbf{h}_u^{(t-1)} \\
&= [\text{diag}(\tilde{\mathbf{a}}^{1T}, \tilde{\mathbf{a}}^{2T}, \dots, \tilde{\mathbf{a}}^{KT}) \parallel \text{diag}(\tilde{\mathbf{a}}'^{1T}, \tilde{\mathbf{a}}'^{2T}, \dots, \tilde{\mathbf{a}}'^{KT})] [\mathbf{W} \mathbf{h}_v^{(t-1)} \parallel \mathbf{W} \mathbf{h}_u^{(t-1)}].
\end{aligned}$$

Therefore, multi-head GAT is an implementation of ExpandingConv as follows:

$$\begin{cases} \alpha_{vu} = \text{softmax}(\text{LeakyReLU}([\text{diag}(\tilde{\mathbf{a}}^{1T}, \tilde{\mathbf{a}}^{2T}, \dots, \tilde{\mathbf{a}}^{KT}) \\ \parallel \text{diag}(\tilde{\mathbf{a}}'^{1T}, \tilde{\mathbf{a}}'^{2T}, \dots, \tilde{\mathbf{a}}'^{KT})][\mathbf{W} \mathbf{h}_v^{(t-1)} \parallel \mathbf{W} \mathbf{h}_u^{(t-1)}])) \\ \mathbf{h}_v^{(t)} = \sigma \left( \frac{1}{K} \mathbf{W} \sum_{j \in \mathcal{N}_v} \text{vec}(\alpha_{vu} \mathbf{h}_u^{(t-1)T}) \right). \end{cases}$$

□

## G COMPARISONS WITH MULTI-AGGREGATOR IMPLEMENTATIONS

ExpandingConv can also be considered as a kind of multi-aggregator scheme. In Equation 5, each row of  $\mathbf{M}_{v_i}$  can be viewed as a weighted aggregator where the weight coefficients are learned from data. Proposition 1 shows that to obtain higher distinguishing strength by utilizing more aggregators, the weight coefficients of newly added aggregators should be linearly independent to all existing aggregators. The distinguishing strength of weighted aggregators is incomparable with basic aggregators. However, since each row of  $\mathbf{M}_{v_i}$  is equivalent to an independent aggregator, one can simply modify the implementation of  $\mathbf{f}_{\text{local}}(u, v)$  to obtain the variant whose distinguishing strength is strict stronger than basic aggregators as follows:

$$\begin{aligned}
& \text{ExpandingConv} \Big|_{\mathbf{m}'_{uv} = [\mathbf{m}_{uv}^{(t)}]_{|1|}} > \text{SUM}, \\
& \text{ExpandingConv} \Big|_{\mathbf{m}'_{uv} = [\mathbf{m}_{uv}^{(t)}]_{|1|} \parallel \frac{1}{|\mathcal{N}(v)|}} > \text{SUM} \otimes \text{MEAN}.
\end{aligned}$$

Compared with leveraging multiple basic aggregators in (Corso et al., 2020) and (Dehmamy et al., 2019), leveraging weighted aggregator allows for variable numbers of aggregators. Meanwhile, the weighted coefficients are learned from data, which can better capture relevant structural patterns.



## H DETAILS OF EXPERIMENTAL SETUP

**Datasets.** Benchmark datasets for graph kernels provided by TU (Kersting et al., 2016) suffer from their small scales of data, making them not sufficient to evaluate the performance of models (Dwivedi et al., 2020). Our evaluations are conducted on graph property predictions datasets ogbg-ppa, ogbg-code, ogbg-molhiv in OGB (Weihua Hu, 2020) and QM9 (Ramakrishnan et al., 2014; Wu et al., 2018; Ruddigkeit et al., 2012) which are large-scale graph datasets including graph classification and graph regression tasks. ogbg-ppa is extracted from the protein-protein association networks with large and densely connected graphs. ogbg-code is a collection of Abstract Syntax Trees (ASTs) obtained from Python method definitions with large and sparse graphs. ogbg-molhiv is molecular property prediction datasets with relative small graphs. QM9 consists 134K small organic molecules with the task to predict 12 targets for each molecule. All data is obtained from pytorch-geometric library (Fey & Lenssen, 2019).

Table 5: Hyperparameter settings for OGB.

	ogbg-ppa		ogbg-molhiv		ogbg-molpcba		ogbg-code	
	ExpC*-1, ExpC-s	CombC*, CombC	ExpC*-1, ExpC-s	CombC*, CombC	ExpC*-1, ExpC-s	CombC*, CombC	ExpC*-1, ExpC-s	CombC*, CombC
batch size	32	32	64	64	128	128	64	64
layers	4	4	3	3	5	5	4	4
hidden	256	256	64	64	512	512	512	512
lr	0.0005	0.0002	0.0001	0.0001	0.0001	0.0001	0.0001	0.0001
step size	20	20	5	5	10	10	5	5
lr decay	0.8	0.7	0.7	0.7	0.6	0.6	0.6	0.6
dropout	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
readout	SUM	SUM	MEAN	MEAN	MEAN	MEAN	MEAN	MEAN

Table 6: Hyperparameter settings for TU.

	COLLAB		REDDIT-BINARY		REDDIT-MULTI-12K	
	ExpC-s	CombC	ExpC-s	CombC	ExpC-s	CombC
batch size	32	32	64	64	64	64
layers	3	3	3	3	3	3
hidden	180	180	256	256	256	256
lr	0.001	0.001	0.001	0.001	0.001	0.001
step size	10	10	10	10	10	10
lr decay	0.8	0.8	0.8	0.8	0.8	0.8
dropout	0.5	0.5	0.5	0.5	0.5	0.5
readout	SUM	SUM	SUM	SUM	SUM	SUM

The shared hyperparameter settings of ExpC\*-1, ExpC-s, CombC\* and CombC on all 12 targets of QM9: batch sizes = 64; lr = 0.0001; step size = 30; lr decay = 0.85; readout = SUM. hidden = 256 for ExpC\*-1 and ExpC-s; hidden = 512 for CombC\* and CombC. Table 7 gives the individual hyperparameter settings of each model on each target, including the number of layers.

Table 7: Number of layers for QM9.

	$\mu$	$\alpha$	$\epsilon_{homo}$	$\epsilon_{lumo}$	$\Delta\epsilon$	$\langle R^2 \rangle$	ZPVE	$U_0$	U	H	G	$C_v$
ExpC*-1,ExpC-s	5	4	5	4	4	4	4	5	4	4	4	4
CombC*,CombC	5	4	5	4	4	5	4	5	4	4	4	4

## I MORE EXPERIMENTAL RESULTS

We present more results of ablation studies on OGB and QM9, which demonstrate the effectiveness of ExpandingConv, CombConv and Re-SUM.

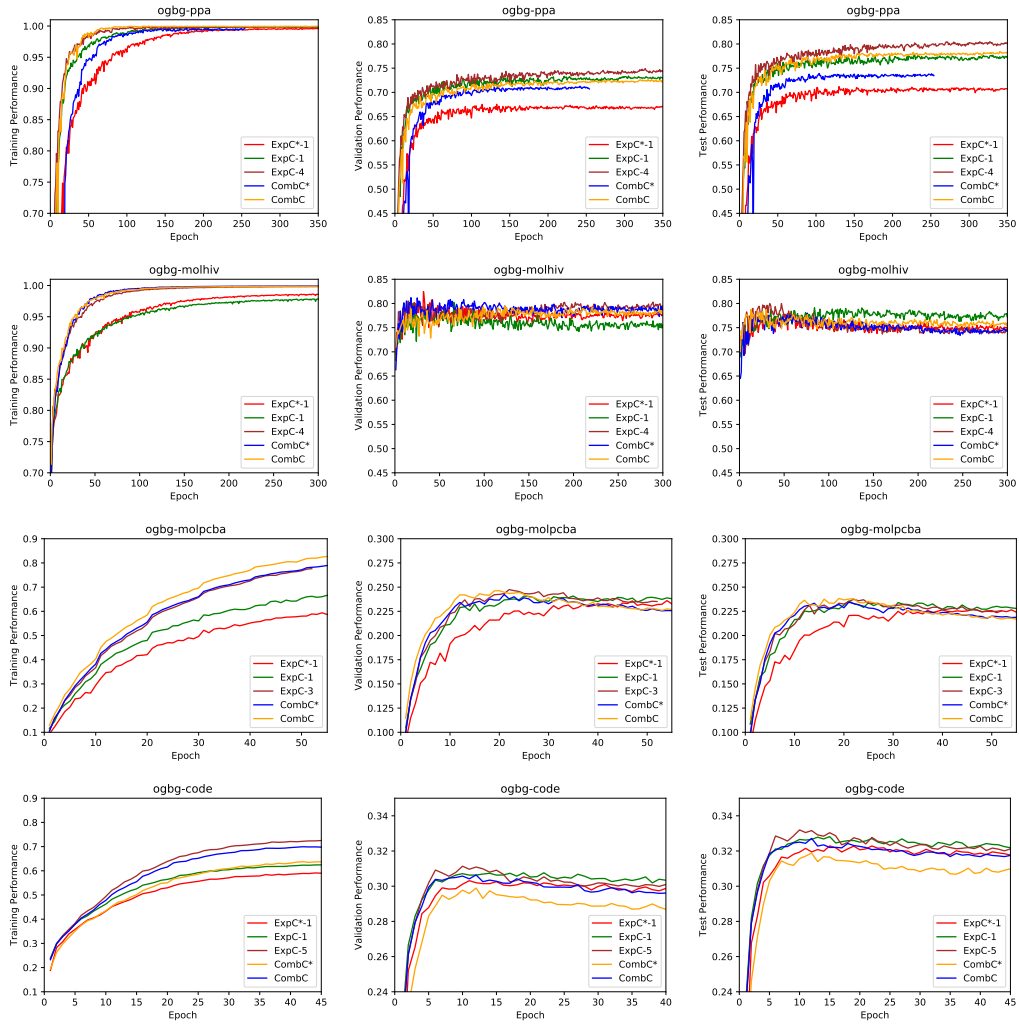


Figure 1: Learning curves on ogbg-ppa, ogbg-molhiv, ogbg-molpcba and ogbg-code.

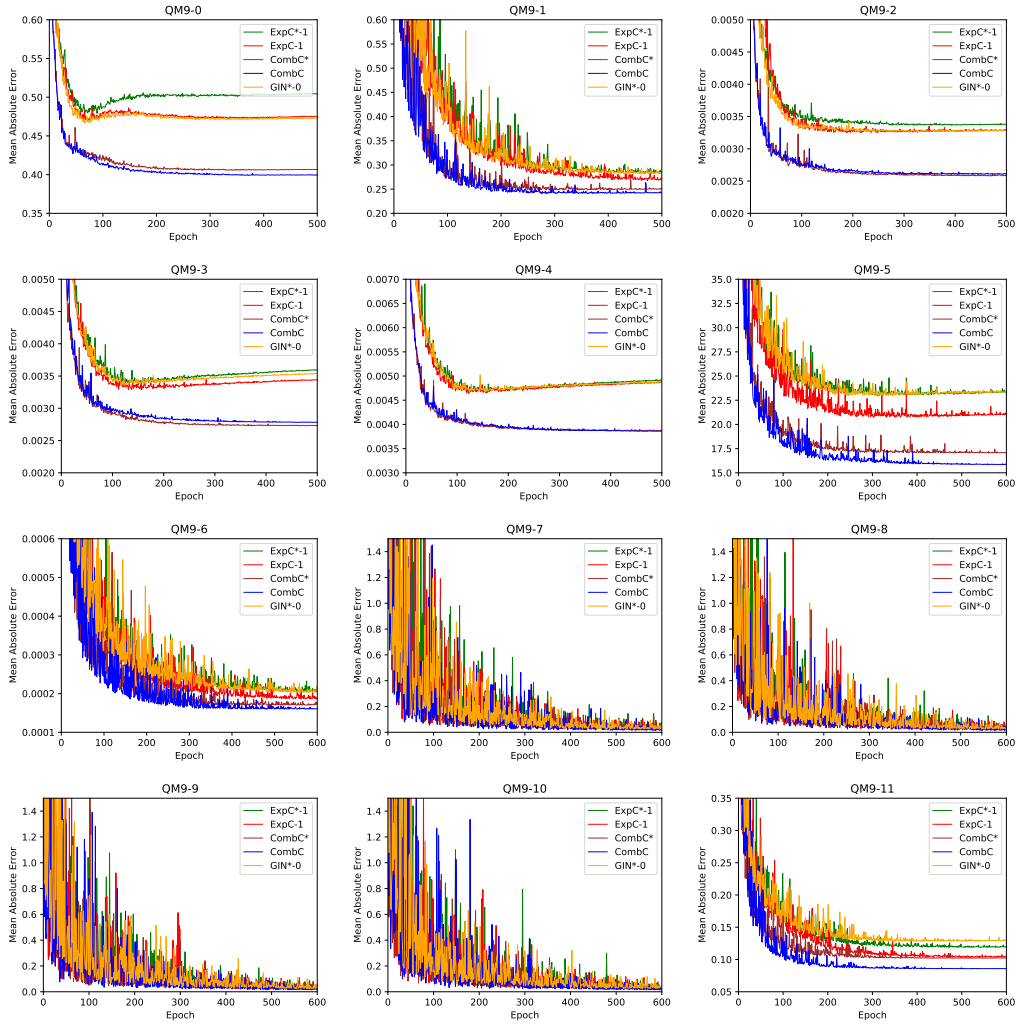


Figure 2: Effectiveness of *Re-SUM* on QM9.

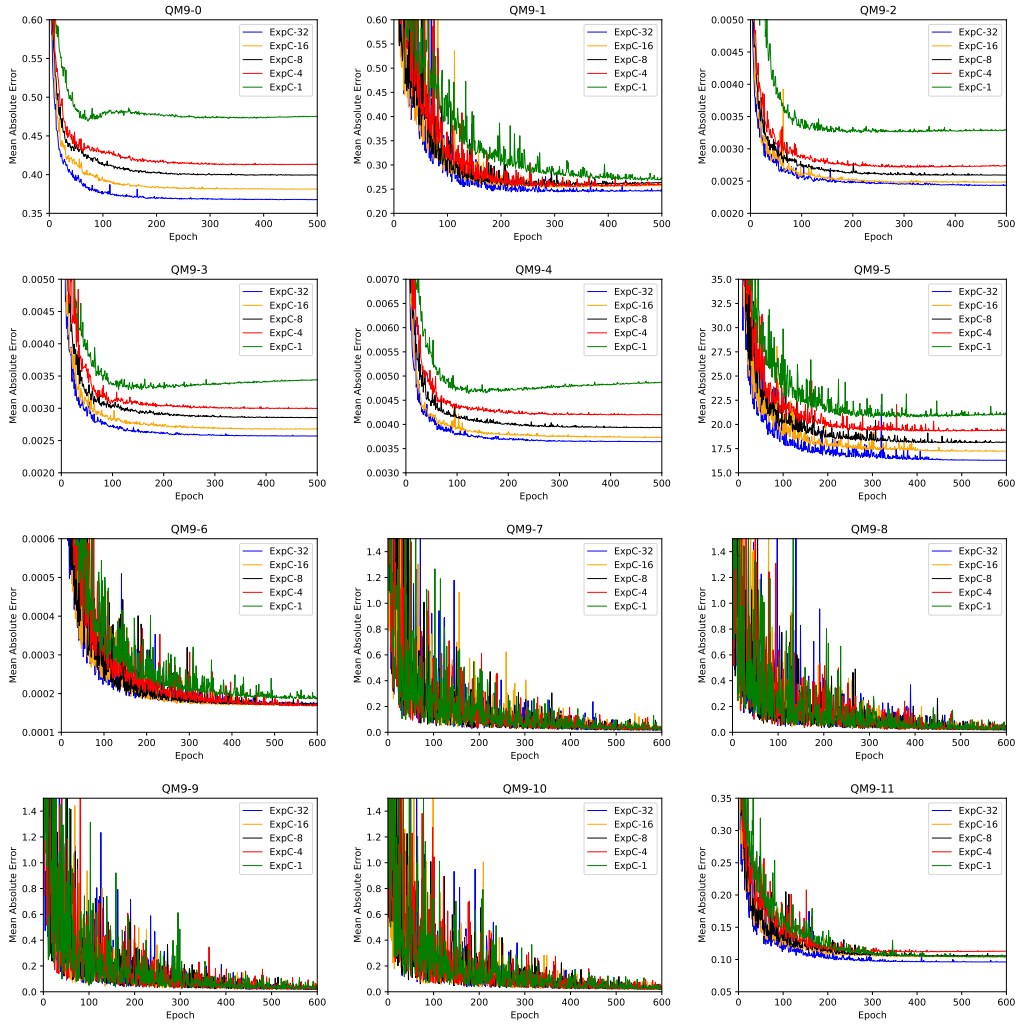


Figure 3: Effectiveness of powerful aggregators on QM9.