

# SPEQNETS: SPARSITY-AWARE PERMUTATION-EQUIVARIANT GRAPH NETWORKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

While graph neural networks have clear limitations in approximating permutation-equivariant functions over graphs, more expressive, higher-order graph neural networks do not scale to large graphs. By introducing new heuristics for the graph isomorphism problem, we devise a class of universal, permutation-equivariant graph networks, which offers a fine-grained control between expressivity and scalability and adapt to the sparsity of the graph. These architectures lead to vastly reduced computation times compared to standard higher-order graph networks while significantly improving over standard graph neural network and graph kernel architectures in terms of predictive performance.

## 1 INTRODUCTION

In recent years, numerous approaches have been proposed for machine learning with graphs—most notably, approaches based on *graph kernels* (Borgwardt et al., 2020; Kriege et al., 2020) or using *graph neural networks* (GNNs) (Chami et al., 2020; Gilmer et al., 2017; Grohe, 2021; Morris et al., 2021). Here, graph kernels based on the *1-dimensional Weisfeiler–Leman algorithm* (1-WL) (Weisfeiler & Leman, 1968), a simple heuristic for the graph isomorphism problem, and corresponding GNNs (Morris et al., 2019; Xu et al., 2019) have recently advanced the state-of-the-art in supervised node- and graph-level learning. However, the 1-WL operates via simple neighborhood aggregation, and the purely local nature of the related approaches misses important patterns in the given data. A provably more powerful algorithm for graph isomorphism testing is the *k-dimensional Weisfeiler–Leman algorithm* (*k*-WL) (Babai, 1979; Cai et al., 1992). The algorithm captures more global, higher-order patterns by iteratively computing a coloring or discrete labeling for *k*-tuples defined over the set of nodes of a given graph based on an appropriately defined notion of adjacency between tuples. See (Kiefer, 2020a;b) for a survey and more background. However, since the algorithm considers all  $n^k$  many *k*-tuples of an *n*-node graph, it does not scale to large, real-world graphs. Moreover, it fixes the cardinality of the neighborhood to  $k \cdot n$ . Hence, a potential *sparsity* of the input graph does not reduce the running time. Further, new neural architectures that possess the same power as the *k*-WL in terms of separating non-isomorphic graphs (Geerts, 2020; Maron et al., 2019b; Morris et al., 2020b) suffer from the same drawbacks.

**Present work** To address the described drawbacks, we introduce a new set of heuristics for the graph isomorphism problem, denoted  $(k, s)$ -LWL, which only considers a subset of all *k*-tuples, namely those *inducing subgraphs with at most s connected components*. We study the effect of *k* and *s* on the expressive power of the heuristics, see Figure 1 for an overview of these theoretical results. Building on these combinatorial insights, we derive corresponding provably expressive, permutation-equivariant neural architectures, denoted  $(k, s)$ -SpeqNets, which offer a more fine-grained trade-off between scalability and expressivity compared to previous architectures based on the *k*-WL. Empirically, we show how our architectures offer vastly reduced computation times while beating baseline GNNs and other higher-order graph networks in terms of predictive performance. See Appendix A for a discussion of related work.

## 2 PRELIMINARIES

We briefly describe the Weisfeiler–Leman algorithm and, along the way, introduce our notation, see Appendix B for details. We let  $[n] := \{1, \dots, n\} \subset \mathbb{N}$  for  $n \geq 1$ , and use  $\{\{ \dots \}\}$  to denote multisets. We also use standard concepts from graph theory (such as graphs, directed graphs, neighbors, trees, and so on). The vertex and the edge set of a graph *G* are denoted by  $V(G)$  and  $E(G)$ , respectively. The

*neighborhood* of  $v$  in  $V(G)$  is  $\delta(v) = \{u \in V(G) \mid (v, u) \in E(G)\}$ . We say that two graphs  $G$  and  $H$  are *isomorphic* ( $G \simeq H$ ) if there exists a bijection  $\varphi: V(G) \rightarrow V(H)$  preserving the adjacency relation, i.e.,  $(u, v)$  is in  $E(G)$  if and only if  $(\varphi(u), \varphi(v))$  is in  $E(H)$ , call  $\varphi$  an *isomorphism* from  $G$  to  $H$ . If the graphs have vertex or edges labels, the isomorphism is additionally required to match these labels. Let  $\mathbf{v}$  be a *tuple* in  $V(G)^k$  for  $k > 0$ , then  $G[\mathbf{v}]$  is the subgraph induced by the elements of  $\mathbf{v}$ , where the nodes are labeled with integers from  $\{1, \dots, k\}$  corresponding to their positions in  $\mathbf{v}$ . A *connected component* of a graph  $G$  is an inclusion-wise maximal subgraph of  $G$  in which every two nodes are connected by paths.

## 2.1 NODE-REFINEMENT ALGORITHMS

In the following, we review the Weisfeiler–Leman algorithm and related variants (Morris et al., 2020b). Let  $k$  be a fixed positive integer and let  $V(G)^k$  denote the set of  $k$ -tuples of nodes of the graph  $G$ . A *coloring* of  $V(G)^k$  is a mapping  $C: V(G)^k \rightarrow \mathbb{N}$ , i.e., we assign a number (color) to every tuple in  $V(G)^k$ . The *initial coloring*  $C_0$  of  $V(G)^k$  is specified by the atomic types of the tuples, i.e., two tuples  $\mathbf{v}$  and  $\mathbf{w}$  in  $V(G)^k$  have the same initial color iff the mapping  $v_i \mapsto w_i$  induces an isomorphism between the labeled subgraphs  $G[\mathbf{v}]$  and  $G[\mathbf{w}]$ . A *color class* corresponding to a color  $c$  is the set of all tuples colored  $c$ , i.e., the set  $C^{-1}(c)$ . For  $j$  in  $[k]$  and  $w$  in  $V(G)$ , let  $\phi_j(\mathbf{v}, w)$  be the  $k$ -tuple obtained by replacing the  $j$ th component of  $\mathbf{v}$  with the node  $w$ . If  $\mathbf{w} = \phi_j(\mathbf{v}, w)$  for some  $w$  in  $V(G)$ , call  $\mathbf{w}$  a  *$j$ -neighbor* of  $\mathbf{v}$ . The *neighborhood* of  $\mathbf{v}$  is the set of all  $\mathbf{w}$  such that  $\mathbf{w} = \phi_j(\mathbf{v}, w)$  holds for some  $j$  in  $[k]$  and a  $w$  in  $V(G)$ .

The *refinement* of a coloring  $C: V(G)^k \rightarrow \mathbb{N}$ , denoted by  $\widehat{C}$ , is a coloring  $\widehat{C}: V(G)^k \rightarrow \mathbb{N}$  defined as follows. For each  $j$  in  $[k]$ , collect the colors of the  $j$ -neighbors of  $\mathbf{v}$  in a multiset  $S_j = \{\{C(\phi_j(\mathbf{v}, w)) \mid w \in V(G)\}\}$ . Then, for a tuple  $\mathbf{v}$ , define  $\widehat{C}(\mathbf{v}) := (C(\mathbf{v}), M(\mathbf{v}))$ , where  $M(\mathbf{v})$  is the  $k$ -tuple  $(S_1, \dots, S_k)$ . For consistency, the strings  $\widehat{C}(\mathbf{v})$  thus obtained are lexicographically sorted and renamed as fresh integers, i.e., ones that have not been used in previous iterations.

**$k$ -dimensional Weisfeiler–Leman** For  $k \geq 2$ , the  $k$ -WL computes a coloring  $C_\infty: V(G)^k \rightarrow \mathbb{N}$  of a given graph  $G$ , as follows.<sup>1</sup> To begin with, the initial coloring  $C_0$  is computed. Then, starting with  $C_0$ , successive refinements  $C_{i+1} = \widehat{C}_i$  are computed until convergence. That is,  $C_{i+1}(\mathbf{v}) = (C_i(\mathbf{v}), M_i(\mathbf{v}))$ , where

$$M_i(\mathbf{v}) := (\{\{C_i(\phi_1(\mathbf{v}, w)) \mid w \in V(G)\}\}, \dots, \{\{C_i(\phi_k(\mathbf{v}, w)) \mid w \in V(G)\}\}).$$

Since the color classes form a partition of  $V(G)^k$ , there must exist a finite  $\ell \leq |V(G)|^k$  such that  $C_\ell = \widehat{C}_\ell$ , i.e., the partition induced by  $C_\ell$  is not refined further.

**Local  $\delta$ - $k$ -dimensional Weisfeiler–Leman algorithm** Morris et al. (2020b) introduced a more efficient modification of the  $k$ -WL, namely the *local  $\delta$ - $k$ -dimensional Weisfeiler–Leman algorithm* ( $\delta$ - $k$ -LWL). In contrast to the  $k$ -WL, the  $\delta$ - $k$ -LWL considers only a subset of the entire neighborhood of a node tuple. Let the tuple  $\mathbf{w} = \phi_j(\mathbf{v}, w)$  be a  $j$ -neighbor of  $\mathbf{v}$ . We say that  $\mathbf{w}$  is a *local  $j$ -neighbor* of  $\mathbf{v}$  if  $w$  is adjacent to the replaced node  $v_j$ . Formally, the  $\delta$ - $k$ -LWL algorithm refines a coloring  $C_i^{k,\delta}$ , obtained after  $i$  rounds of  $\delta$ - $k$ -LWL, via,

$$M_i^\delta(\mathbf{v}) := (\{\{C_i^{k,\delta}(\phi_1(\mathbf{v}, w)) \mid w \in \delta(v_1)\}\}, \dots, \{\{C_i^{k,\delta}(\phi_k(\mathbf{v}, w)) \mid w \in \delta(v_k)\}\}), \quad (1)$$

hence considering only the local  $j$ -neighbors of the tuple  $\mathbf{v}$  in each iteration. The coloring function for the  $\delta$ - $k$ -LWL is then defined by

$$C_{i+1}^{k,\delta}(\mathbf{v}) := (C_i^{k,\delta}(\mathbf{v}), M_i^\delta(\mathbf{v})). \quad (2)$$

We define the 1-WL to be the  $\delta$ -1-LWL. Morris et al. (2020b) also defined the  $\delta$ - $k$ -LWL<sup>+</sup> and showed that the  $\delta$ - $k$ -LWL<sup>+</sup> is slightly more powerful than the  $k$ -WL in distinguishing non-isomorphic graphs, see Appendix B for details.

**Comparing  $k$ -WL variants** Let  $A_1$  and  $A_2$  denote two node-refinement algorithms. We write  $A_1 \sqsubseteq A_2$  if  $A_1$  distinguishes between all non-isomorphic pairs that  $A_2$  distinguishes, and  $A_1 \equiv A_2$  if both  $A_1 \sqsubseteq A_2$  and  $A_2 \sqsubseteq A_1$  hold. The corresponding strict relation is denoted by  $\sqsubset$ .

**Kernels based on node-refinement algorithms** After running the  $k$ -WL (and the other node-refinement algorithms), the concatenation of the histogram of colors in each iteration can be used as a feature vector in a kernel computation. Specifically, in the histogram, for every color  $c$  in  $\mathbb{N}$ , an entry contains the number of nodes or  $k$ -tuples colored  $c$ .

<sup>1</sup>We define the 1-WL in the next subsection.

### 3 THE $(k, s)$ -LWL ALGORITHM AND SPEQNETS

Since both  $k$ -WL and its local variant  $\delta$ - $k$ -LWL consider all  $k$ -tuples of a graph, they do not scale to large graphs for growing  $k$ . To address this issue, we introduce the  $(k, s)$ -LWL, it considers a subset of all  $k$ -tuples, namely those inducing subgraphs with at most  $s$  *connected components*. Formally, let  $G$  be a graph. Then  $\#\text{com}(G)$  denotes the number of (connected) components of  $G$ . Further, let  $k \geq 1$  and  $s$  in  $[k]$ , then

$$V(G)_s^k := \{\mathbf{v} \in V(G)^k \mid \#\text{com}(G[\mathbf{v}]) \leq s\}$$

is the set of  $(k, s)$ -tuples of nodes, i.e,  $k$ -tuples which induce (sub-)graphs with at most  $s$  (connected) components. In contrast to the algorithms of Section 2.1, the  $(k, s)$ -LWL colors tuples from  $V(G)_s^k$  instead of the entire  $V(G)^k$ . Just as in Section 2.1, a coloring of  $V(G)_s^k$  is a mapping  $C_i^{k,s}: V(G)_s^k \rightarrow \mathbb{N}$  for  $i \geq 0$ , assigning a number (color) to every tuple in  $V(G)_s^k$ , and Equation (1) is replaced with

$$M_i^{\delta,k,s}(\mathbf{v}) := (\{\{C_i^{k,s}(\phi_1(\mathbf{v}, w)) \mid w \in \delta(v_1); \phi_1(\mathbf{v}, w) \in V(G)_s^k\}\}, \dots, \{\{C_i^{k,s}(\phi_k(\mathbf{v}, w)) \mid w \in \delta(v_k); \phi_k(\mathbf{v}, w) \in V(G)_s^k\}\}).$$

The following results show that the  $(k, 1)$ -LWL forms a *hierarchy*, i.e., the algorithm becomes more expressive as  $k$  increases.

**Theorem 1.** For  $k \geq 1$ , it holds that

$$(k+1, 1)\text{-LWL} \sqsubset (k, 1)\text{-LWL}.$$

Moreover, we also show that the  $(k, 2)$ -LWL is more expressive than the  $(k, 1)$ -LWL.

**Proposition 2.** For  $k \geq 2$ , it holds that

$$(k, 2)\text{-LWL} \sqsubset (k, 1)\text{-LWL}.$$

Further, the following theorem yields that increasing the parameter  $s$  results in higher expressivity.

**Theorem 3.** For  $k \geq 2$ , it holds that

$$(k, k)\text{-LWL} \sqsubset (k, 2)\text{-LWL}.$$

See Appendix C.2 for an analysis of the running time of the  $(k, s)$ -LWL, showing that its running time on an bounded-degree,  $n$ -vertex graph is  $\tilde{O}(n^s)$  instead of the usual  $\tilde{O}(n^k)$  for the  $k$ -WL, for fixed  $k$  and  $s$ .

#### 3.1 SPEQNETS

We can now leverage the above combinatorial insights to derive sparsity-aware, permutation-equivariant graph networks, denoted  $(k, s)$ -SpeqNets. Given a node-labeled graph  $G$ , let each  $(k, s)$ -tuple  $\mathbf{v}$  in  $V(G)_s^k$  be annotated with an initial feature  $f^{(0)}(\mathbf{v})$  determined by its (labeled) isomorphism type, e.g., a one-hot encoding. In each layer  $t > 0$ , we compute a new feature  $f^{(t)}(\mathbf{v})$  as

$$f_{\text{mrg}}^{W_1} \left( f^{(t-1)}(\mathbf{v}), f_{\text{agg}}^{W_2} \left( \left( \{ \{ f^{(t-1)}(\phi_i(\mathbf{v}, w)) \mid w \in \delta(v_i) \text{ and } \phi_i(\mathbf{v}, w) \in V(G)_s^k \} \}_{i \in [k]} \right) \right) \right), \quad (3)$$

in  $\mathbb{R}^{1 \times e}$ , where  $W_1^{(t)}$  and  $W_2^{(t)}$  are learnable parameter matrices from  $\mathbb{R}^{d \times e}$ .<sup>2</sup> Here,  $f_{\text{mrg}}^{W_2}$  and  $f_{\text{agg}}^{W_1}$  are differentiable functions to merge and aggregate the relevant feature information, respectively. See Appendix D for further theoretical properties of the architecture.

## 4 EXPERIMENTAL EVALUATION

Here, we aim to investigate the learning performance of the  $(k, s)$ -LWL-based kernel and neural architectures. Concretely, we aim to answer the following questions.

**Q1** Do the  $(k, s)$ -LWL-based algorithms lead to improved classification and regression scores on graph-level benchmark datasets compared standard baselines?

**Q2** How does the  $(k, s)$ -SpeqNets compare to standard GNN baselines for node classification?

Table 1: Classification accuracies in percent and standard deviations.

Method	Dataset								
	ENZYMES	IMDB-BINARY	IMDB-MULTI	MUTAG	NCII	PROTEINS	PTC_MR	REDDIT-BINARY	
Baseline	GR	29.9 ±0.8	59.3 ±0.9	39.2 ±0.6	72.5 ±1.7	66.2 ±0.2	71.5 ±0.5	56.6 ±1.3	59.7 ±0.5
	SP	40.3 ±0.9	58.7 ±0.6	39.7 ±0.3	81.7 ±1.5	74.1 ±0.2	75.8 ±0.7	59.6 ±1.5	84.5 ±0.2
	1-WL	50.6 ±1.2	72.5 ±0.8	50.0 ±0.8	75.9 ±2.0	84.4 ±0.3	73.1 ±0.6	59.3 ±2.1	73.4 ±0.9
	WLOA	57.1 ±0.8	73.2 ±0.4	49.8 ±0.4	83.4 ±1.2	85.2 ±0.2	73.0 ±0.9	60.3 ±1.9	88.3 ±0.4
GNN	Gin- $\epsilon$	38.7 ±1.5	72.9 ±0.7	49.7 ±0.7	84.1 ±1.4	77.7 ±0.8	72.2 ±0.6	55.2 ±1.7	89.8 ±0.4
	Gin- $\epsilon$ -JK	39.3 ±1.6	73.0 ±1.1	49.6 ±0.7	83.4 ±2.0	78.3 ±0.3	72.2 ±0.7	56.0 1.3±	90.4 ±0.4
$k$ -WL	2-WL	37.0 ±1.0	68.1 ±1.7	47.5 ±0.7	85.7 ±1.6	66.9 ±0.3	75.2 ±0.4	60.5 ±1.1	OOM
	3-WL	42.3 ±1.1	67.1 ±1.5	46.8 ±0.8	85.4 ±1.5	OOT	OOT	59.0 ±2.0	OOM
	$\delta$ -2-LWL	55.9 ±1.0	73.0 ±0.7	50.1 ±0.9	85.6 ±1.4	84.6 ±0.3	75.1 ±0.5	61.7 ±2.4	89.4 ±0.6
	$\delta$ -2-LWL <sup>+</sup>	53.9 ±1.4	75.6 ±1.0	62.7 ±1.4	84.1 ±2.1	<b>91.3</b> ±0.3	79.2 ±1.2	61.6 ±1.3	91.4 ±0.4
	$\delta$ -3-LWL	<b>58.2</b> ±1.2	72.6 ±0.9	49.0 ±1.2	84.1 ±1.6	83.2 ±0.4	OOM	60.7 ±2.2	OOM
	$\delta$ -3-LWL <sup>+</sup>	56.5 ±1.4	76.1 ±1.2	64.3 ±1.2	85.4 ±1.8	82.7 ±0.4	OOM	61.5 ±1.8	OOM
$(k, s)$ -LWL	(2, 1)-LWL	53.7 ±1.7	73.5 ±0.8	50.8 ±0.7	84.2 ±1.7	82.8 ±0.3	73.2 ±0.6	55.9 ±2.4	76.9 ±0.6
	(2, 1)-LWL <sup>+</sup>	51.6 ±1.8	73.7 ±1.1	55.4 ±0.9	79.6 ±3.4	81.9 ±0.3	76.0 ±0.9	60.2 ±2.1	<b>94.7</b> ±0.3
	(3, 1)-LWL	53.4 ±1.4	74.6 ±1.0	51.3 ±0.6	85.3 ±2.4	81.4 ±0.5	72.9 ±1.1	60.2 ±1.7	OOM
	(3, 1)-LWL <sup>+</sup>	57.0 ±1.9	<b>87.1</b> ±0.6	<b>67.1</b> ±1.1	79.2 ±1.5	89.8 ±0.4	<b>81.2</b> ±0.8	59.2 ±2.0	OOM
	(3, 2)-LWL	56.4 ±0.7	73.5 ±0.5	49.7 ±0.6	<b>86.4</b> ±2.6	84.9 ±0.4	75.1 ±0.9	61.9 ±2.4	OOM
	(3, 2)-LWL <sup>+</sup>	55.8 ±1.7	78.1 ±1.4	59.5 ±1.0	84.5 ±1.9	89.4 ±0.3	78.8 ±0.6	<b>62.3</b> ±3.3	OOM

Table 2: Additional experimental results for graph regression and node classification.

(a) Mean MAE (mean std. MAE, logMAE) on large-scale (multi-target) molecular regression tasks.

Method	Dataset	
	ALCHEMY (10K)	QM9
GINE- $\epsilon$	0.180 ±0.006 -1.958 ±0.047	0.079 ±0.003 -3.430 ±0.080
(2, 1)-SpeqNet	0.169 ±0.005 -2.010 ±0.056	0.078 ±0.007 -2.947 ±0.171
(2, 2)-SpeqNet	<b>0.115</b> ±0.001 -2.722 ±0.054	<b>0.029</b> ±0.001 -4.081 ±0.058
(3, 1)-SpeqNet	0.180 ±0.011 -1.914 ±0.097	0.068 ±0.003 -3.397 ±0.086
(3, 2)-SpeqNet	<b>0.115</b> ±0.002 -2.767 ±0.079	OOT

(b) Classification accuracies and standard deviations for node classification.

Method	Dataset		
	CORNELL	TEXAS	WISCONSIN
GCN	56.5 ±0.9	58.2 ±0.8	50.9 ±0.7
SDRF + Undirected	57.5 ±0.3	<b>70.4</b> ±0.6	61.6 ±0.9
(2, 1)-SpeqNet	63.9 ±1.7	66.8 ±0.9	67.7 ±2.2
(2, 2)-SpeqNet	<b>67.9</b> ±1.7	67.3 ±2.0	<b>68.4</b> ±2.2
(3, 1)-SpeqNet	61.8 ±3.3	68.3 ±1.3	60.4 ±2.8

**Q3** To what extent does the  $(k, s)$ -LWL reduce computation times compared to  $k$ -WL-based architectures? See Appendix E for details on the experiments.

**Results and discussion** In the following, we answer questions **Q1** to **Q3**.

**A1 Kernels** See Table 1. The  $(k, s)$ -LWL for  $k, s$  in  $\{2, 3\}$  significantly improves the classification accuracy compared to the  $k$ -WL and the  $\delta$ - $k$ -WL, and the other kernel baselines, while being on par with or better than the  $\delta$ -2-LWL and  $\delta$ -3-LWL. The  $(k, s)$ -LWL and  $(k, s)$ -LWL<sup>+</sup> achieve a new state-of-the-art on five out of eight datasets. Our algorithms also perform vastly better than the neural baselines.

**Neural architectures** See Table 2. On both datasets, all  $(k, s)$ -SpeqNet architectures beat the GNN baseline. On the ALCHEMY dataset the (2, 2)-SpeqNet and (3, 1)-SpeqNet perform best, while on the QM9 dataset the (2, 2)-SpeqNet performs best by a large margin.

**A2** See Table 2. Over all three datasets, the  $(k, s)$ -SpeqNet architectures improve over the GCN baseline. Specifically, over all datasets, the (2, 1)-SpeqNet and the (2, 2)-SpeqNet lead to an increase of at least 7% in accuracy. For example, both architectures beat the GCN baseline by at least 17% on the WISCONSIN dataset. Further, the  $(k, s)$ -SpeqNet architectures lead to better accuracies compared to the SDRF architecture.

**A3** See Appendix F.

## 5 CONCLUSION

To circumvent the exponential running time requirements of  $k$ -WL, we introduced a new heuristic for the graph isomorphism problem, namely the  $(k, s)$ -LWL. By varying the parameters  $k$  and  $s$ , the  $(k, s)$ -LWL offers a tradeoff between scalability and expressivity and, unlike the  $k$ -WL, takes into account the potential graph’s sparsity. Based on these combinatorial insights, we designed provably expressive machine-learning architectures. Empirically, we showed that such architectures lead to state-of-the-art results in node- and graph-level classification regimes while obtaining promising results on graph-level regression tasks.

<sup>2</sup>For clarity of presentation, we omit biases.

## REFERENCES

- R. Abboud, Í. Í. Ceylan, M. Grohe, and T. Lukaszewicz. The surprising power of graph neural networks with random node initialization. *ArXiv preprint*, 2020.
- B. M. Anderson, T.-S. Hy, and R. Kondor. Cormorant: Covariant molecular neural networks. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 14510–14519, 2019.
- V. Arvind, J. Köbler, G. Rattan, and O. Verbitsky. On the power of color refinement. In *International Symposium on Fundamentals of Computation Theory*, pp. 339–350, 2015.
- V. Arvind, F. Fuhlbrück, J. Köbler, and O. Verbitsky. On Weisfeiler-Leman invariance: Subgraph counts and related graph properties. In *International Symposium on Fundamentals of Computation Theory*, pp. 111–125, 2019.
- A. Atserias and E. N. Maneva. Sherali-adams relaxations and indistinguishability in counting logics. *SIAM Journal on Computing*, (1):112–137, 2013.
- A. Atserias, L. Mancinska, D. E. Roberson, R. Sámal, S. Severini, and A. Varvitsiotis. Quantum and non-signalling graph isomorphisms. *Journal of Combinatorial Theory, Series B*, pp. 289–328, 2019.
- W. Azizian and M. Lelarge. Characterizing the expressive power of invariant and equivariant graph neural networks. *ArXiv preprint*, 2020.
- L. Babai. Lectures on graph isomorphism. University of Toronto, Department of Computer Science. Mimeographed lecture notes, October 1979, 1979.
- L. Babai. Graph isomorphism in quasipolynomial time. In *ACM SIGACT Symposium on Theory of Computing*, pp. 684–697, 2016.
- M. Balcilar, P. Héroux, B. Gaüzère, P. Vasseur, S. Adam, and P. Honeine. Breaking the limits of message passing graph neural networks. In *International Conference on Machine Learning*, pp. 599–608, 2021.
- P. Barceló, E. V. Kostylev, M. Monet, J. Pérez, J. L. Reutter, and J. Pablo Silva. The logical expressiveness of graph neural networks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- P. Barceló, F. Geerts, J. L. Reutter, and M. Ryschkov. Graph neural networks with local graph parameters. *ArXiv preprint*, 2021.
- I. I. Baskin, V. A. Palyulin, and N. S. Zefirov. A neural device for searching direct correlations between structures and properties of chemical compounds. *Journal of Chemical Information and Computer Sciences*, (4):715–721, 1997.
- D. Beaini, S. Passaro, V. Létourneau, W. L. Hamilton, G. Corso, and P. Liò. Directional graph networks. *ArXiv preprint*, 2020.
- C. Berkholz, P. S. Bonsma, and M. Grohe. Tight lower and upper bounds for the complexity of canonical colour refinement. *Theory of Computing Systems*, (4):581–614, 2017.
- B. Bevilacqua, F. Frasca, D. Lim, B. Srinivasan, C. Cai, G. Balamurugan, M. M. Bronstein, and H. Maron. Equivariant subgraph aggregation networks. *ArXiv preprint*, 2021.
- C. Bodnar, F. Frasca, N. Otter, Y. G. Wang, P. Liò, G. Montúfar, and M. M. Bronstein. Weisfeiler and Lehman go cellular: CW networks. *ArXiv preprint*, 2021a.
- C. Bodnar, F. Frasca, Y. Wang, N. Otter, G. F. Montufar, P. Lió, and M. Bronstein. Weisfeiler and Lehman go topological: Message passing simplicial networks. In *International Conference on Machine Learning*, pp. 1026–1037, 2021b.
- K. M. Borgwardt and H.-P. Kriegel. Shortest-path kernels on graphs. In *IEEE International Conference on Data Mining*, pp. 74–81, 2005.

- K. M. Borgwardt, M. Elisabetta Ghisu, F. Llinares-López, L. O’Bray, and B. Rieck. Graph kernels: State-of-the-art and future challenges. *Foundations and Trends in Machine Learning*, (5-6), 2020.
- G. Bouritsas, F. Frasca, S. Zafeiriou, and M. M. Bronstein. Improving graph neural network expressivity via subgraph isomorphism counting. *ArXiv preprint*, 2020.
- J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and deep locally connected networks on graphs. In *International Conference on Learning Representation*, 2014.
- R. G. Busacker and T. L. Saaty. *Finite graphs and networks: an introduction with applications*. 1965.
- J. Cai, M. Fürer, and N. Immerman. An optimal lower bound on the number of variables for graph identifications. *Combinatorica*, (4):389–410, 1992.
- I. Chami, S. Abu-El-Haija, B. Perozzi, C. Ré, and K. Murphy. Machine learning on graphs: A model and comprehensive taxonomy. *ArXiv preprint*, 2020.
- C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, pp. 27:1–27:27, 2011. ACM.
- G. Chen, Pengfei Chen, C.-Y. Hsieh, C.-K. Lee, B. Liao, R. Liao, W. Liu, J. Qiu, Q. Sun, J. Tang, R. S. Zemel, and S. Zhang. Alchemy: A quantum chemistry dataset for benchmarking AI models. *ArXiv preprint*, 2019a.
- Z. Chen, S. Villar, L. Chen, and J. Bruna. On the equivalence between graph isomorphism testing and function approximation with gnns. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 15868–15876, 2019b.
- Z. Chen, L. Chen, S. Villar, and J. Bruna. Can graph neural networks count substructures? In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- L. Cotta, C. Morris, and B. Ribeiro. Reconstruction for powerful graph representations. In *Advances in Neural Information Processing Systems*, 2021.
- G. Dasoulas, L. Dos Santos, K. Scaman, and A. Virmaux. Coloring graph neural networks for node disambiguation. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pp. 2126–2132, 2020.
- M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 3837–3845, 2016.
- H. Dell, M. Grohe, and G. Rattan. Lovász meets Weisfeiler and Leman. In *International Colloquium on Automata, Languages, and Programming*, pp. 40:1–40:14, 2018.
- D. Duvenaud, D. Maclaurin, J. Aguilera-Iparraguirre, R. Gómez-Bombarelli, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp. 2224–2232, 2015.
- M. Fey and J. E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *International Conference on Learning Representations, Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- M. Fürer. On the combinatorial power of the Weisfeiler-Lehman algorithm. In *International Conference on Algorithms and Complexity*, pp. 260–271, 2017.
- T. Gärtner, P. A. Flach, and S. Wrobel. On graph kernels: Hardness results and efficient alternatives. In *Computational Learning Theory and Kernel Machines, 16th Annual Conference on Computational Learning Theory and 7th Kernel Workshop, COLT/Kernel 2003, Washington, DC, USA, August 24-27, 2003, Proceedings*, pp. 129–143, 2003.

- F. Geerts. The expressive power of  $k$ th-order invariant graph networks. *ArXiv preprint*, 2020.
- F. Geerts, F. Mazowiecki, and G. A. Pérez. Let’s agree to degree: Comparing graph convolutional networks in the message-passing framework. *ArXiv preprint*, 2020.
- J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pp. 1263–1272, 2017.
- M. Grohe. *Descriptive Complexity, Canonisation, and Definable Graph Structure Theory*. 2017.
- M. Grohe. Word2vec, Node2vec, Graph2vec, X2vec: Towards a theory of vector embeddings of structured data. *ArXiv preprint*, 2020.
- M. Grohe. The logic of graph neural networks. In *ACM/IEEE Symposium on Logic in Computer Science*, pp. 1–17, 2021.
- M. Grohe and M. Otto. Pebble games and linear equations. *Journal of Symbolic Logic*, (3):797–844, 2015.
- M. Grohe, P. Schweitzer, and Wiebking D. Deep Weisfeiler Leman. *ArXiv preprint*, 2020.
- W. L. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 1024–1034, 2017.
- M. Heimann, T. Safavi, and D. Koutra. Distribution of node embeddings as multiresolution features for graphs. In *IEEE International Conference on Data Mining*, pp. 289–298, 2019.
- N. Immerman and E. Lander. Describing graphs: A first-order approach to graph canonization. In *Complexity Theory Retrospective: In Honor of Juris Hartmanis on the Occasion of His Sixtieth Birthday, July 5, 1988*, pp. 59–81, 1990.
- F. D. Johansson and D. P. Dubhashi. Learning with similarity functions on graphs using matchings of geometric embeddings. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, pp. 467–476, 2015.
- H. Kashima, K. Tsuda, and A. Inokuchi. Marginalized kernels between labeled graphs. In *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003), August 21-24, 2003, Washington, DC, USA*, pp. 321–328, 2003.
- S. Kiefer. *Power and Limits of the Weisfeiler-Leman Algorithm*. PhD thesis, Department of Computer Science, RWTH Aachen University, 2020a.
- S. Kiefer. The Weisfeiler-Leman algorithm: an exploration of its power. *ACM SIGLOG News*, (3):5–27, 2020b.
- S. Kiefer and B. D. McKay. The iteration number of colour refinement. In *International Colloquium on Automata, Languages, and Programming*, pp. 73:1–73:19, 2020.
- S. Kiefer and P. Schweitzer. Upper bounds on the quantifier depth for graph differentiation in first order logic. In *ACM/IEEE Symposium on Logic in Computer Science*, pp. 287–296, 2016.
- S. Kiefer, P. Schweitzer, and E. Selman. Graphs identified by logics with counting. In *International Symposium on Mathematical Foundations of Computer Science*, pp. 319–330, 2015.
- S. Kiefer, I. Ponomarenko, and P. Schweitzer. The Weisfeiler-Leman dimension of planar graphs is at most 3. *Journal of the ACM*, (6):44:1–44:31, 2019.
- T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- D. B. Kireev. Chemnet: A novel neural network based method for graph/property mapping. *Journal of Chemical Information and Computer Sciences*, (2):175–180, 1995. ACS.

- J. Klicpera, J. Groß, and S. Günnemann. Directional message passing for molecular graphs. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- R. Kondor and H. Pan. The multiscale Laplacian graph kernel. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 2982–2990, 2016.
- N. M. Kriege, P.-L. Giscard, and R. C. Wilson. On valid optimal assignment kernels and applications to graph classification. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 1615–1623, 2016.
- N. M. Kriege, M. Neumann, C. Morris, K. Kersting, and P. Mutzel. A unifying view of explicit and implicit feature maps for structured data: Systematic studies of graph kernels. *ArXiv preprint*, 2017.
- N. M. Kriege, C. Morris, A. Rey, and C. Sohler. A property testing framework for the theoretical expressivity of graph kernels. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pp. 2348–2354, 2018.
- N. M. Kriege, F. D. Johansson, and C. Morris. A survey on graph kernels. *Applied Network Science*, (1):6, 2020.
- P. Li, Y. Wang, H. Wang, and J. Leskovec. Distance encoding: Design provably more powerful neural networks for graph representation learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- M. Lichter, I. Ponomarenko, and P. Schweitzer. Walk refinement, walk logic, and the iteration number of the Weisfeiler-Leman algorithm. In *ACM/IEEE Symposium on Logic in Computer Science*, pp. 1–13, 2019.
- T. Maehara and H. NT. A simple proof of the universality of invariant/equivariant graph neural networks. *ArXiv preprint*, 2019.
- P. N. Malkin. Sherali–adams relaxations of graph isomorphism polytopes. *Discrete Optimization*, pp. 73–97, 2014.
- H. Maron, H. Ben-Hamu, H. Serviansky, and Y. Lipman. Provably powerful graph networks. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 2153–2164, 2019a.
- H. Maron, E. Fetaya, N. Segol, and Y. Lipman. On the universality of invariant networks. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pp. 4363–4371, 2019b.
- C. Merkwirth and T. Lengauer. Automatic generation of complementary descriptors with molecular graph networks. *Journal of Chemical Information and Modeling*, (5):1159–1168, 2005.
- A. Micheli. Neural network for graphs: A contextual constructive approach. *IEEE Transactions on Neural Networks*, (3):498–511, 2009.
- A. Micheli and A. S. Sestito. A new neural network model for contextual processing of graphs. In *Italian Workshop on Neural Nets Neural Nets and International Workshop on Natural and Artificial Immune Systems*, pp. 10–17, 2005.
- F. Monti, D. Boscaini, J. Masci, E. Rodolà, J. Svoboda, and M. M. Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 5425–5434, 2017.
- C. Morris, K. Kersting, and P. Mutzel. Glocalised Weisfeiler-Lehman kernels: Global-local feature maps of graphs. In *IEEE International Conference on Data Mining*, pp. 327–336, 2017.



- C. Morris, M. Ritzert, M. Fey, W. L. Hamilton, J. E. Lenssen, G. Rattan, and M. Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pp. 4602–4609, 2019.
- C. Morris, N. M. Kriege, F. Bause, K. Kersting, P. Mutzel, and M. Neumann. TUDataset: A collection of benchmark datasets for learning with graphs. *ArXiv preprint*, 2020a.
- C. Morris, G. Rattan, and P. Mutzel. Weisfeiler and Leman go sparse: Towards higher-order graph embeddings. In *Advances in Neural Information Processing Systems*, 2020b.
- C. Morris, Y. L., H. Maron, B. Rieck, N. M. Kriege, M. Grohe, M. Fey, and K. Borgwardt. Weisfeiler and Leman go machine learning: The story so far. *ArXiv preprint*, 2021.
- R. L. Murphy, B. Srinivasan, V. A. Rao, and B. Ribeiro. Relational pooling for graph representations. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pp. 4663–4673, 2019.
- H. Nguyen and T. Maehara. Graph homomorphism convolution. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, pp. 7306–7316, 2020.
- G. Nikolentzos, P. Meladianos, and M. Vazirgiannis. Matching node embeddings for graph similarity. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, pp. 2429–2435, 2017.
- G. Nikolentzos, P. Meladianos, S. Limnios, and M. Vazirgiannis. A degeneracy framework for graph similarity. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pp. 2595–2601, 2018.
- P. A. Papp, L. Faber K. Martinkus, and R. Wattenhofer. DropGNN: Random dropouts increase the expressiveness of graph neural networks. In *Advances in Neural Information Processing Systems*, 2021.
- H. Pei, B. Wei, K. Chen-Chuan Chang, Y. Lei, and B. Yang. Geom-gcn: Geometric graph convolutional networks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- R. Ramakrishnan, O. Dral, P., M. Rupp, and O. Anatole von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data*, 2014. Nature.
- B. Rieck, C. Bock, and Karsten M. Borgwardt. A persistent weisfeiler-lehman procedure for graph classification. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pp. 5448–5458, 2019.
- R. Sato, M. Yamada, and H. Kashima. Random features strengthen graph neural networks. *ArXiv preprint*, 2020.
- F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, (1):61–80, 2009.
- N. Shervashidze, S. V. N. Vishwanathan, T. H. Petri, K. Mehlhorn, and K. M. Borgwardt. Efficient graphlet kernels for large graph comparison. In *International Conference on Artificial Intelligence and Statistics*, pp. 488–495, 2009.
- N. Shervashidze, P. Schweitzer, E. J. van Leeuwen, K. Mehlhorn, and K. M. Borgwardt. Weisfeiler-Lehman graph kernels. *Journal of Machine Learning Research*, pp. 2539–2561, 2011.
- A. Sperduti and A. Starita. Supervised neural networks for the classification of structures. *IEEE Transactions on Neural Networks*, (2):714–35, 1997. IEEE.
- R. Talak, S. Hu, L. Peng, and L. Carlone. Neural trees for learning on graphs. *ArXiv preprint*, 2021.
- E. H. Thiede, W. Zhou, and R. Kondor. Autobahn: Automorphism-based graph neural nets. *ArXiv preprint*, 2021.

- M. Togninalli, M. E. Ghisu, F. Llinares-López, B. Rieck, and K. M. Borgwardt. Wasserstein weisfeiler-lehman graph kernels. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 6436–6446, 2019.
- J. Tönshoff, M. Ritzert, H. Wolf, and M. Grohe. Graph learning with 1D convolutions on random walks. *ArXiv preprint*, 2021.
- J. Topping, F. Di Giovanni, B. P. Chamberlain, X. Dong, and M. M. Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature. *ArXiv preprint*, 2021.
- G. Valiente. *Algorithms on Trees and Graphs*. 2002.
- P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph attention networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- S. Verma and Z.-L. Zhang. Hunt for the unique, stable, sparse and fast feature learning on graphs. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 88–98, 2017.
- C. Vignac, A. Loukas, and P. Frossard. Building powerful and equivariant graph neural networks with structural message-passing. In *Advances in Neural Information Processing Systems*, 2020.
- O. Vinyals, S. Bengio, and M. Kudlur. Order matters: Sequence to sequence for sets. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- E. Wagstaff, F. B. Fuchs, M. Engelcke, M. A. Osborne, and I. Posner. Universal approximation of functions on sets. *ArXiv preprint*, 2021.
- B. Weisfeiler. *On Construction and Identification of Graphs*. 1976.
- B. Weisfeiler and A. Leman. The reduction of a graph to canonical form and the algebra which appears therein. *Nauchno-Technicheskaya Informatsia*, 2(9):12–16, 1968. English translation by G. Ryabov is available at [https://www.itl.zcu.cz/wl2018/pdf/wl\\_paper\\_translation.pdf](https://www.itl.zcu.cz/wl2018/pdf/wl_paper_translation.pdf).
- Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, and V. Pande. MoleculeNet: A benchmark for molecular machine learning. *Chemical Science*, pp. 513–530, 2018.
- K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019.
- P. Yanardag and S. V. N. Vishwanathan. A structural smoothing framework for robust graph comparison. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp. 2134–2142, 2015a.
- P. Yanardag and S. V. N. Vishwanathan. Deep graph kernels. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, pp. 1365–1374, 2015b.
- J. You, J. Gomes-Selman, R. Ying, and J. Leskovec. Identity-aware graph neural networks. *ArXiv preprint*, 2021.
- M. Zhang and P. Li. Nested graph neural networks. *ArXiv preprint*, 2021.
- L. Zhao, W. Jin, L. Akoglu, and N. Shah. From stars to subgraphs: Uplifting any GNN with local structure awareness. *ArXiv preprint*, 2021.

## A RELATED WORK

In the following, we review related work from graph kernels, GNNs and graph theory.

**Graph kernels** Historically, kernel methods—which implicitly or explicitly map graphs to elements of a Hilbert space—have been the dominant approach for supervised learning on graphs. Important early work in this area includes random-walk based kernels (Gärtner et al., 2003; Kashima et al., 2003; Kriege et al., 2017) and kernels based on shortest paths (Borgwardt & Kriegel, 2005). More recently, developments in the field have emphasized scalability, focusing on techniques that bypass expensive Gram matrix computations by using explicit feature maps, see, e.g., (Shervashidze et al., 2011). Morris et al. (2017) devised a local, set-based variant of the  $k$ -WL and a corresponding kernel. However, the approach is (provably) weaker than the tuple-based algorithm. Further, Morris et al. (2020a) proposed kernels based on the  $\delta$ - $k$ -LWL.

Yanardag & Vishwanathan (2015a) successfully employed Graphlet (Shervashidze et al., 2009), and Weisfeiler–Leman kernels within frameworks for smoothed (Yanardag & Vishwanathan, 2015a) and deep graph kernels (Yanardag & Vishwanathan, 2015b). Other recent work focuses on assignment-based (Johansson & Dubhashi, 2015; Kriege et al., 2016; Nikolentzos et al., 2017), spectral (Kondor & Pan, 2016; Verma & Zhang, 2017), graph decomposition (Nikolentzos et al., 2018), randomized binning approaches (Heimann et al., 2019), and the extension of kernels based on the 1-WL (Rieck et al., 2019; Togninalli et al., 2019). For a theoretical investigation of graph kernels, see (Kriege et al., 2018), and for a thorough survey of graph kernels, see (Borgwardt et al., 2020; Kriege et al., 2020).

**GNNs** Recently, GNNs (Gilmer et al., 2017; Scarselli et al., 2009) emerged as an alternative to graph kernels. Notable instances of this architecture include, e.g., (Duvenaud et al., 2015; Hamilton et al., 2017; Velickovic et al., 2018), which can be subsumed under the message-passing framework introduced in (Gilmer et al., 2017). Also, approaches based on spectral information were introduced in, e.g., (Defferrard et al., 2016; Bruna et al., 2014; Kipf & Welling, 2017; Monti et al., 2017)—all of which descend from early work in (Kireev, 1995; Baskin et al., 1997; Micheli & Sestito, 2005; Merkwirth & Lengauer, 2005; Micheli, 2009; Sperduti & Starita, 1997; Scarselli et al., 2009).

**Limits of GNNs and more expressive architectures** Recently, connections of GNNs to Weisfeiler–Leman type algorithms have been shown (Azizian & Lelarge, 2020; Barceló et al., 2020; Chen et al., 2019b; Geerts et al., 2020; Geerts, 2020; Maehara & NT, 2019; Maron et al., 2019a; Morris et al., 2019; Xu et al., 2019). Specifically, (Morris et al., 2019; Xu et al., 2019) showed that the expressive power of any possible GNN architecture is limited by the 1-WL in terms of distinguishing non-isomorphic graphs.

Triggered by the above results, a large set of papers proposed architectures to overcome the expressivity limitations of the 1-WL. Morris et al. (2019) introduced  $k$ -dimensional GNNs ( $k$ -GNN) which rely on a message-passing scheme between subgraphs of cardinality  $k$ . Similar to (Morris et al., 2017), the paper employed a local, set-based (neural) variant of the  $k$ -WL. Later, this was refined in (Maron et al., 2019a; Azizian & Lelarge, 2020) by introducing  $k$ -order folklore graph neural networks ( $k$ -FGNN), which are equivalent to the folklore or oblivious variant of the  $k$ -WL (Grohe, 2021; Morris et al., 2021) in terms of distinguishing non-isomorphic graphs. Subsequently, Morris et al. (2020b) introduced neural architectures based on the  $\delta$ - $k$ -LWL, which only considers a subset of the neighborhood from the  $k$ -WL, taking sparsity of the underlying graph (to some extent) into account. Although more scalable, the algorithm reaches computational exhaustion on large-scale graphs since it considers all  $n^k$  tuples of size  $k$ . Chen et al. (2019b) connected the theory of universal approximations of permutation-invariant functions and the graph isomorphism viewpoint and introduced a variation of the 2-WL. See (Morris et al., 2021) for an in-depth survey on this topic.

Recent works have extended the expressive power of GNNs, e.g., by encoding node identifiers (Murphy et al., 2019; Vignac et al., 2020), leveraging random features (Abboud et al., 2020; Dasoulas et al., 2020; Sato et al., 2020), subgraph information (Bevilacqua et al., 2021; Bouritsas et al., 2020; Cotta et al., 2021; Papp et al., 2021; Thiede et al., 2021; You et al., 2021; Zhang & Li, 2021; Zhao et al., 2021), homomorphism counts (Barceló et al., 2021; Nguyen & Maehara, 2020), spectral information (Balcilar et al., 2021), simplicial and cellular complexes (Bodnar et al., 2021b;a), random walks (Tönshoff et al., 2021), graph decompositions (Talak et al., 2021), distance (Li et al., 2020) and directional information (Beaini et al., 2020).

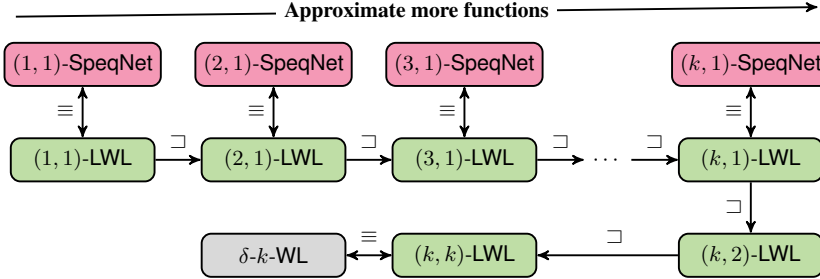


Figure 1: Overview of the power of the proposed algorithms and neural architectures. The green and red nodes represent algorithms proposed in the present work. Forward arrows point to more powerful algorithms or neural architectures.

$A \sqsubset B$  ( $A \equiv B$ ): algorithm  $A$  is strictly more powerful (equally powerful) than  $B$ .

However, up to a few exceptions, all of the above approaches only overcome limitations of the 1-WL or the 2-WL, and do not induce a hierarchy of provably powerful, permutation-equivariant neural architectures aligned with the  $k$ -WL hierarchy.

**Theory** The Weisfeiler–Leman algorithm constitutes one of the earliest and most natural approaches to isomorphism testing (Weisfeiler, 1976; Weisfeiler & Leman, 1968), having been heavily investigated by the theory community over the last few decades (Grohe, 2017). Moreover, the fundamental nature of the  $k$ -WL is evident from a variety of connections to other fields such as logic, optimization, counting complexity, and quantum computing. The power and limitations of  $k$ -WL can be neatly characterized in terms of logic and descriptive complexity (Babai, 1979; Immerman & Lander, 1990), Sherali-Adams relaxations of the natural integer linear program for the graph isomorphism problem (Atserias & Maneva, 2013; Grohe & Otto, 2015; Malkin, 2014), homomorphism counts (Dell et al., 2018), and quantum isomorphism games (Atserias et al., 2019). In their seminal paper, Cai et al. (1992) showed that, for each  $k$ , there exists a pair of non-isomorphic graphs of size  $\mathcal{O}(k)$  that are not distinguished by the  $k$ -WL. (Kiefer, 2020a;b) gives a thorough survey of these results. For  $k = 1$ , the power of the algorithm has been completely characterized (Arvind et al., 2015; Kiefer et al., 2015). Moreover, upper bounds on the running time (Berkholz et al., 2017) and the number of iterations for  $k = 1$  (Kiefer & McKay, 2020) and for the folklore  $k = 2$  (Kiefer & Schweitzer, 2016; Lichter et al., 2019) have been shown. For  $k$  in  $\{1, 2\}$ , Arvind et al. (2019) studied the abilities of the (folklore)  $k$ -WL to detect and count fixed subgraphs, extending the work of Fürer (2017). The former was refined in (Chen et al., 2020). Kiefer et al. (2019) showed that the folklore 3-WL completely captures the structure of planar graphs. The algorithm (for logarithmic  $k$ ) plays a prominent role in the recent result of Babai (2016) improving the best-known running time for the graph isomorphism problem. Recently, Grohe et al. (2020) introduced the framework of Deep Weisfeiler–Leman algorithms, which allow the design of a more powerful graph isomorphism test than Weisfeiler–Leman type algorithms. Finally, the emerging connections between the Weisfeiler–Leman paradigm and graph learning are described in two recent surveys (Grohe, 2020; Morris et al., 2021).

## B PRELIMINARIES

As usual, let  $[n] := \{1, \dots, n\} \subset \mathbb{N}$  for  $n \geq 1$ , and we use  $\{\!\{ \dots \}\!\}$  to denote multisets, i.e., the generalization of sets allowing for multiple instances for each of its elements.

**Graphs** A graph  $G$  is a pair  $(V(G), E(G))$  with finite sets of nodes  $V(G)$  and edges  $E(G) \subseteq \{\{u, v\} \subseteq V \mid u \neq v\}$ . If not otherwise stated, we set  $n := |V(G)|$ . For ease of notation, we denote the edge  $\{u, v\}$  in  $E(G)$  by  $(u, v)$  or  $(v, u)$ . In the case of *directed graphs*,  $E \subseteq \{(u, v) \in V \times V \mid u \neq v\}$ . A *labeled graph*  $G$  is a triple  $(V, E, \ell)$  with a label function  $\ell: V(G) \cup E(G) \rightarrow \mathbb{N}$ . Then  $\ell(x)$  is a *label* of  $x$  for  $x$  in  $V(G) \cup E(G)$ . The *neighborhood* of  $v$  in  $V(G)$  is denoted by  $\delta(v) = \{u \in V(G) \mid (v, u) \in E(G)\}$  and the *degree* of a node  $v$  is  $|\delta(v)|$ . Let  $S \subseteq V(G)$  then  $G[S] = (S, E_S)$  is the *subgraph induced by  $S$* , where  $E_S = \{(u, v) \in E(G) \mid u, v \in S\}$ . A *connected component* of a graph  $G$  is an inclusion-wise maximal subgraph of  $G$  in which every two nodes are connected by paths. A *tree* is a connected graph without cycles. A *rooted tree* is an oriented tree with a designated node called *root*, in which the edges are directed away from the root. Let  $p$  be a node in a rooted tree. Then we call its out-neighbors *children* with *parent*  $p$ . We denote an undirected *cycle* on  $k$  nodes by  $C_k$ . Given two graphs  $G$  and  $H$  with disjoint node sets, we denote their disjoint union by  $G \dot{\cup} H$ .

Two graphs  $G$  and  $H$  are *isomorphic* and we write  $G \simeq H$  if there exists a bijection  $\varphi: V(G) \rightarrow V(H)$  preserving the adjacency relation, i.e.,  $(u, v)$  is in  $E(G)$  if and only if  $(\varphi(u), \varphi(v))$  is in  $E(H)$ . Then  $\varphi$  is an *isomorphism* between  $G$  and  $H$ . Moreover, we call the equivalence classes induced by  $\simeq$  *isomorphism types*, and denote the isomorphism type of  $G$  by  $\tau_G$ . In the case of labeled graphs, we additionally require that  $\ell(v) = \ell(\varphi(v))$  for  $v$  in  $V(G)$  and  $\ell((u, v)) = \ell((\varphi(u), \varphi(v)))$  for  $(u, v)$  in  $E(G)$ . Let  $\mathbf{v}$  be a *tuple* in  $V(G)^k$  for  $k > 0$ , then  $G[\mathbf{v}]$  is the subgraph induced by the multiset of elements of  $\mathbf{v}$ , where the nodes are labeled with integers from  $\{1, \dots, k\}$  corresponding to their positions in  $\mathbf{v}$ .

**Equivariance** For  $n > 0$ , let  $S_n$  denote the set of all permutations of  $[n]$ , i.e., the set of all bijections from  $[n]$  to itself. For  $\sigma$  in  $S_n$  and a graph  $G$  such that  $V(G) = [n]$ , let  $G_\sigma = \sigma \cdot G$  be the graph such that  $V(\sigma \cdot G) = \{v_{\sigma^{-1}(1)}, \dots, v_{\sigma^{-1}(n)}\}$  and  $E(G_\sigma) = \{(v_{\sigma^{-1}(i)}, v_{\sigma^{-1}(j)}) \mid (v_i, v_j) \in E(G)\}$ . That is, applying the permutation  $\sigma$  reorders the nodes. Hence, for two isomorphic graphs  $G$  and  $H$  on the same vertex set, i.e.,  $G \simeq H$ , there exists  $\sigma$  in  $S_n$  such that  $\sigma \cdot G = H$ .

Let  $\mathcal{G}$  denote the set of all graphs, and let  $\mathcal{G}_n$  denote the set of all graphs on  $n$  nodes. A function  $f: \mathcal{G} \rightarrow \mathbb{R}$  is *invariant* if for every  $n > 0$  and every  $\sigma$  in  $S_n$  and graph  $G$ ,  $f(\sigma \cdot G) = f(G)$ . A function  $f: \mathcal{G} \rightarrow \mathbb{R}$  is *equivariant* if for every  $n > 0$ ,  $f(\mathcal{G}_n) \subseteq \mathbb{R}$  and for every  $\sigma$  in  $S_n$ ,  $f(\sigma \cdot G) = \sigma \cdot f(G)$ .

**Kernels** A *kernel* on a non-empty set  $\mathcal{X}$  is a symmetric, positive semidefinite function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . Equivalently, a function  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a kernel if there is a *feature map*  $\phi: \mathcal{X} \rightarrow \mathcal{H}$  to a Hilbert space  $\mathcal{H}$  with inner product  $\langle \cdot, \cdot \rangle$ , such that  $k(x, y) = \langle \phi(x), \phi(y) \rangle$  for all  $x$  and  $y$  in  $\mathcal{X}$ . A *graph kernel* is a kernel on the set  $\mathcal{G}$  of all graphs.

## B.1 NODE-REFINEMENT ALGORITHMS (EXTENDED)

In the following, we briefly describe the Weisfeiler–Leman algorithm and related variants (Morris et al., 2020b). Let  $k$  be a fixed positive integer. There exist two definitions of the  $k$ -WL, the so-called oblivious  $k$ -WL and folklore or non-oblivious  $k$ -WL, in literature, see, e.g., (Grohe, 2021). There is a subtle difference in how they aggregate neighborhood information. Within the graph learning community, it is customary to abbreviate the oblivious  $k$ -WL as  $k$ -WL, a convention that we follow in this paper.

We proceed to the definition of the  $k$ -WL. Let  $V(G)^k$  denote the set of  $k$ -tuples of nodes of the graph  $G$ . A *coloring* of  $V(G)^k$  is a mapping  $C: V(G)^k \rightarrow \mathbb{N}$ , i.e., we assign a number (color) to every tuple in  $V(G)^k$ . The *initial coloring*  $C_0$  of  $V(G)^k$  is specified by the atomic types of the tuples, i.e., two tuples  $\mathbf{v}$  and  $\mathbf{w}$  in  $V(G)^k$  have the same initial color iff mapping  $v_i \mapsto w_i$  induces an isomorphism between the labeled subgraphs  $G[\mathbf{v}]$  and  $G[\mathbf{w}]$ . Note that, given a tuple  $\mathbf{v}$  in  $V(G)^k$ , we can upper-bound the running time of the computation of this initial coloring for  $\mathbf{v}$  by  $\mathcal{O}(k^2)$ . A *color class* corresponding to a color  $c$  is the set of all tuples colored  $c$ , i.e., the set  $C^{-1}(c)$ .

For  $j$  in  $[k]$  and  $w$  in  $V(G)$ , let  $\phi_j(\mathbf{v}, w)$  be the  $k$ -tuple obtained by replacing the  $j$ th component of  $\mathbf{v}$  with the node  $w$ . That is,  $\phi_j(\mathbf{v}, w) = (v_1, \dots, v_{j-1}, w, v_{j+1}, \dots, v_k)$ . If  $\mathbf{w} = \phi_j(\mathbf{v}, w)$  for some  $w$  in  $V(G)$ , call  $\mathbf{w}$  a  $j$ -neighbor of  $\mathbf{v}$ . The *neighborhood* of  $\mathbf{v}$  is the set of all  $\mathbf{w}$  such that  $\mathbf{w} = \phi_j(\mathbf{v}, w)$  for some  $j$  in  $[k]$  and a  $w \in V(G)$ .

The *refinement* of a coloring  $C: V(G)^k \rightarrow \mathbb{N}$ , denoted by  $\widehat{C}$ , is a coloring  $\widehat{C}: V(G)^k \rightarrow \mathbb{N}$  defined as follows. For each  $j$  in  $[k]$ , collect the colors of the  $j$ -neighbors of  $\mathbf{v}$  in a multiset  $S_j = \{C(\phi_j(\mathbf{v}, w)) \mid w \in V(G)\}$ . Then, for a tuple  $\mathbf{v}$ , define

$$\widehat{C}(\mathbf{v}) := (C(\mathbf{v}), M(\mathbf{v})),$$

where  $M(\mathbf{v})$  is the  $k$ -tuple  $(S_1, \dots, S_k)$ . For consistency, the strings  $\widehat{C}(\mathbf{v})$  thus obtained are lexicographically sorted and renamed as integers, not used in previous iterations. Observe that the new color  $\widehat{C}(\mathbf{v})$  of  $\mathbf{v}$  is solely dictated by the color histogram of the neighborhood of  $\mathbf{v}$ . In general, a different mapping  $M(\cdot)$  could be used, depending on the neighborhood information that we would like to aggregate. We will refer to a mapping  $M(\cdot)$  as an *aggregation map*.

**$k$ -dimensional Weisfeiler–Leman** For  $k \geq 2$ , the  $k$ -WL computes a coloring  $C_\infty: V(G)^k \rightarrow \mathbb{N}$  of a given graph  $G$ , as follows.<sup>3</sup> To begin with, the initial coloring  $C_0$  is computed. Then, starting with  $C_0$ , successive refinements  $C_{i+1} = \widehat{C}_i$  are computed until convergence. That is,

$$C_{i+1}(\mathbf{v}) = (C_i(\mathbf{v}), M_i(\mathbf{v})),$$

<sup>3</sup>We define the 1-WL in the next subsection.

where

$$M_i(\mathbf{v}) = (\{\{C_i(\phi_1(\mathbf{v}, w)) \mid w \in V(G)\}, \dots, \{\{C_i(\phi_k(\mathbf{v}, w)) \mid w \in V(G)\}\}). \quad (4)$$

The successive refinement steps are also called *rounds* or *iterations*. Since the disjoint union of the color classes form a partition of  $V(G)^k$ , there must exist a finite  $\ell \leq |V(G)|^k$  such that  $C_\ell = \widehat{C}_\ell$ , i.e., the partition induced by  $C_\ell$  cannot be refined further. The  $k$ -WL outputs  $C_\ell$  as the *stable coloring*  $C_\infty$ .

The  $k$ -WL *distinguishes* two graphs  $G$  and  $H$  if, upon running the  $k$ -WL on their disjoint union  $G \dot{\cup} H$ , there exists a color  $c$  in  $\mathbb{N}$  in the stable coloring such that the corresponding color class  $S_c$  satisfies

$$|V(G)^k \cap S_c| \neq |V(H)^k \cap S_c|,$$

i.e., there the numbers of  $c$ -colored tuples in  $V(G)^k$  and  $V(H)^k$  differ. Two graphs that are distinguished by the  $k$ -WL must be non-isomorphic, because the algorithm is defined in an isomorphism-invariant way.

Finally, the application of different aggregation maps  $M(\cdot)$  yield related versions of  $k$ -WL. For example, setting  $M(\cdot)$  to be

$$M^F(\mathbf{v}) = (\{C(\phi_1(\mathbf{v}, w)), \dots, C(\phi_k(\mathbf{v}, w)) \mid w \in V(G)\},$$

yields the so-called folklore-version of  $k$ -WL (see e.g. (Cai et al., 1992)). It is known that the oblivious version of the  $k$ -WL is as powerful as the folklore version of the  $(k-1)$ -WL (Grohe, 2021).

### Local $\delta$ - $k$ -dimensional Weisfeiler–Leman algorithm

Morris et al. (2020b) introduced a more efficient variant of the  $k$ -WL, namely the *local  $\delta$ - $k$ -dimensional Weisfeiler–Leman algorithm* ( $\delta$ - $k$ -LWL). In contrast to the  $k$ -WL, the  $\delta$ - $k$ -LWL considers only a subset of the entire neighborhood of a node tuple. Let the tuple  $\mathbf{w} = \phi_j(\mathbf{v}, w)$  be a  $j$ -neighbor of  $\mathbf{v}$ . We say that  $\mathbf{w}$  is a *local  $j$ -neighbor* of  $\mathbf{v}$  if  $w$  is adjacent to the replaced node  $v_j$ . Otherwise, the tuple  $\mathbf{w}$  is a *global  $j$ -neighbor* of  $\mathbf{v}$ . The  $\delta$ - $k$ -LWL considers only local neighbors during the neighborhood aggregation process, and discards any information about the global neighbors. Formally, the  $\delta$ - $k$ -LWL algorithm refines a coloring  $C_i^{k,\delta}$  (obtained after  $i$  rounds of  $\delta$ - $k$ -LWL) via the aggregation function,

$$M_i^\delta(\mathbf{v}) = (\{\{C_i^{k,\delta}(\phi_1(\mathbf{v}, w)) \mid w \in \delta(v_1)\}, \dots, \{\{C_i^{k,\delta}(\phi_k(\mathbf{v}, w)) \mid w \in \delta(v_k)\}\}), \quad (5)$$

hence considering only the local  $j$ -neighbors of the tuple  $\mathbf{v}$  in each iteration. The coloring function for the  $\delta$ - $k$ -LWL is then defined by

$$C_{i+1}^{k,\delta}(\mathbf{v}) = (C_i^{k,\delta}(\mathbf{v}), M_i^\delta(\mathbf{v})). \quad (6)$$

We define the 1-WL to be the  $\delta$ -1-LWL, which is commonly known as Color Refinement or Naive Node Classification.<sup>4</sup> Hence, we can equivalently define

$$C_{i+1}^{1,\delta}(v) = (C_i^{1,\delta}(v), \{\{C_i^{1,\delta}(w) \mid w \in \delta(v)\}\}). \quad (7)$$

for a node  $v$  in  $V(G)$ .

Morris et al. (2020b) also defined the  $\delta$ - $k$ -LWL<sup>+</sup>, a minor variation of the  $\delta$ - $k$ -LWL. Formally, the  $\delta$ - $k$ -LWL<sup>+</sup> refines a coloring  $C_i$  (obtained after  $i$  rounds) via the aggregation function

$$M_i^{\delta,+}(\mathbf{v}) = (\{\{C_i^{k,\delta}(\phi_1(\mathbf{v}, w)), \#_i^1(\mathbf{v}, \phi_1(\mathbf{v}, w)) \mid w \in \delta(v_1)\}, \dots, \{\{C_i^{k,\delta}(\phi_k(\mathbf{v}, w)), \#_i^k(\mathbf{v}, \phi_k(\mathbf{v}, w)) \mid w \in \delta(v_k)\}\}), \quad (8)$$

instead of the  $\delta$ - $k$ -LWL aggregation defined in Equation (5). Here, we set

$$\#_i^j(\mathbf{v}, \mathbf{x}) := |\{\mathbf{w} : \mathbf{w} \sim_j \mathbf{v}, C_i^{k,\delta}(\mathbf{w}) = C_i^{k,\delta}(\mathbf{x})\}|, \quad (9)$$

where  $\mathbf{w} \sim_j \mathbf{v}$  denotes that  $\mathbf{w}$  is a  $j$ -neighbor of  $\mathbf{v}$ , for  $j$  in  $[k]$ . Essentially,  $\#_i^j(\mathbf{v}, \mathbf{x})$  counts the number of  $j$ -neighbors (local or global) of  $\mathbf{v}$  which have the same color as  $\mathbf{x}$  under the coloring  $C_i$  (i.e., after  $i$  rounds). Morris et al. (2020b) showed that the  $\delta$ - $k$ -LWL<sup>+</sup> is slightly more powerful than the  $k$ -WL in distinguishing non-isomorphic graphs.

<sup>4</sup>Strictly speaking, the 1-WL and Color Refinement are two different algorithms. That is, the 1-WL considers neighbors and non-neighbors to update the coloring, resulting in a slightly higher expressivity when distinguishing nodes in a given graph, see (Grohe, 2021) for details. For brevity, we consider both algorithms to be equivalent.

## C THE $(k, s)$ -LWL ALGORITHM (EXTENDED)

Since both  $k$ -WL and its local variant  $\delta$ - $k$ -LWL consider all  $k$ -tuples of a graph, they do not scale to large graphs for larger  $k$ . Specifically, for an  $n$ -node graph, the memory requirement is  $\Omega(n^k)$ . Further, since the  $k$ -WL considers the graph structure only at initialization, it does not adapt to its sparsity, i.e., it does not run faster for sparser graphs. To address this issue, we introduce the  $(k, s)$ -LWL. The algorithm offers more fine-grained control over the trade-off between expressivity and scalability by only considering a subset of all  $k$ -tuples, namely those inducing subgraphs with at most  $s$  *connected components*. This combinatorial algorithm will be the basis of the permutation-equivariant neural architectures, see below.

Formally, let  $G$  be a graph, then  $\#\text{com}(G)$  denotes the number of (connected) components of  $G$ . Further, let  $k \geq 1$  and  $1 \leq s \leq k$ , then

$$V(G)_s^k := \{\mathbf{v} \in V(G)^k \mid \#\text{com}(G[\mathbf{v}]) \leq s\}$$

is the set of  $(k, s)$ -tuples of nodes, i.e.  $k$ -tuples which induce (sub-)graphs with at most  $s$  (connected) components.

In contrast to the algorithms of Section 2.1, the  $(k, s)$ -LWL colors tuples from  $V(G)_s^k$  instead of the entire  $V(G)^k$ . Hence, analogously to Section 2.1, a coloring of  $V(G)_s^k$  is a mapping  $C_i^{k,s}: V(G)_s^k \rightarrow \mathbb{N}$  for  $i \geq 0$ , assigning a number (color) to every tuple in  $V(G)_s^k$ . The initial coloring  $C_0^{k,s}$  of  $V(G)_s^k$  is defined in the same way as before, i.e., specified by the isomorphism types of the tuples, see Section 2.1. Subsequently, the coloring is updated using the  $\delta$ - $k$ -LWL aggregation map, see Equation (1). Hence, the  $(k, s)$ -LWL is a variant of the  $\delta$ - $k$ -LWL considering only  $(k, s)$ -tuples, i.e., Equation (1) is replaced with

$$\begin{aligned} M_i^{\delta,k,s}(\mathbf{v}) := & \left( \{C_i^{k,s}(\phi_1(\mathbf{v}, w)) \mid w \in \delta(v_1) \text{ and } \phi_1(\mathbf{v}, w) \in V(G)_s^k\}, \right. \\ & \left. \dots, \{C_i^{k,s}(\phi_k(\mathbf{v}, w)) \mid w \in \delta(v_k) \text{ and } \phi_k(\mathbf{v}, w) \in V(G)_s^k\} \right), \end{aligned} \quad (10)$$

i.e.,  $M_i^\delta(\mathbf{v})$  restricted to colors of  $(k, s)$ -tuples. The stable coloring  $C_\infty^{k,s}$  is defined analogously to the stable coloring  $C_\infty^k$ . In the following two subsections, we investigate the properties of the algorithm in detail.

Analogously to the  $\delta$ - $k$ -LWL<sup>+</sup>, we also define the  $(k, s)$ -LWL<sup>+</sup> using

$$\begin{aligned} M_i^{\delta,+}(\mathbf{v}) = & \left( \{C_i^{k,s}(\phi_1(\mathbf{v}, w)), \#_{i,s}^1(\mathbf{v}, \phi_1(\mathbf{v}, w)) \mid w \in \delta(v_1) \text{ and } \phi_1(\mathbf{v}, w) \in V(G)_s^k\}, \dots, \right. \\ & \left. \{C_i^{k,s}(\phi_k(\mathbf{v}, w)), \#_{i,s}^k(\mathbf{v}, \phi_k(\mathbf{v}, w)) \mid w \in \delta(v_k) \text{ and } \phi_k(\mathbf{v}, w) \in V(G)_s^k\} \right), \end{aligned}$$

where the function

$$\#_{i,s}^j(\mathbf{v}, \mathbf{x}) = \left| \{\mathbf{w}: \mathbf{w} \sim_j \mathbf{v}, C_i^{k,s}(\mathbf{w}) = C_i^{k,s}(\mathbf{x}) \text{ and } \mathbf{w} \in V(G)_s^k\} \right|,$$

restricts  $\#_{i,s}^j(\mathbf{v}, \mathbf{x})$  to  $(k, s)$ -tuples.

### C.1 EXPRESSIVITY

Here, we investigate the expressivity of the  $(k, s)$ -LWL, i.e., its ability to distinguish non-isomorphic graphs, for different choices of  $k$  and  $s$ . Below we will leverage these results to devise universal, permutation-equivariant graph networks. We start off with the following simple observation. Since the  $(k, k)$ -LWL colors all  $k$ -tuples, it is equal to the  $\delta$ - $k$ -LWL.

**Observation 1.** Let  $k \geq 1$ , then

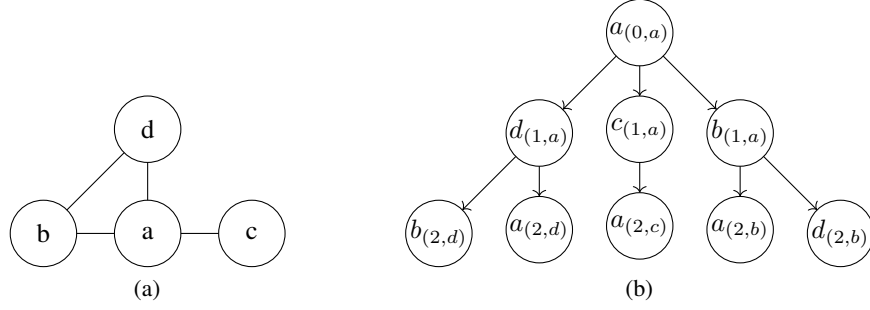
$$(k, k)\text{-LWL} \equiv \delta\text{-}k\text{-LWL}.$$

The following results shows that the  $(k, 1)$ -LWL forms a *hierarchy*, i.e., the algorithm gets more expressive as  $k$  increases.

**Theorem 4.** Let  $k \geq 1$ , then

$$(k+1, 1)\text{-LWL} \sqsubset (k, 1)\text{-LWL}.$$

To prove Theorem 4, we introduce the  $(k, s)$ -tuple graph. It essentially contains the set of all  $(k, s)$ -tuples as nodes, where each node  $v_{\mathbf{t}}$  is labeled by the isomorphism type of the  $(k, s)$ -tuple  $\mathbf{t}$ . We join two nodes by an edge, labeled  $j$ , if the underlying  $(k, s)$ -tuples are  $j$ -neighbors. The formal definition of the  $(k, s)$ -tuple graph is as follows. (Recall that  $\tau$  denotes isomorphism types.)

Figure 2: Illustration of the unrolling operation around the node  $a$  for  $i = 2$ .

**Definition 5.** Let  $G$  be a graph and let  $k \geq 1$ , and  $s$  in  $[k]$ . Further, let  $\mathbf{s}$  and  $\mathbf{t}$  be tuples in  $V(G)_s^k$ . Then the directed, labeled  $(k, s)$ -tuple graph  $T_s^k(G) = (V_T, E_T, \ell_T)$  has node set  $V_T = \{v_t \mid \mathbf{t} \in V(G)_s^k\}$ , and

$$(v_s, v_t) \in E_T \iff \mathbf{t} = \phi_j(\mathbf{s}, w) \text{ holds for some } j \text{ in } [k] \text{ and some } w \text{ in } V(G). \quad (11)$$

We set  $\ell_T((v_s, v_t)) := j$  if  $\mathbf{t}$  is a local  $j$ -neighbor of  $\mathbf{s}$ , and let  $\ell_T(v_s) := \tau_{G[\mathbf{s}]}$ .

Given a graph  $G$  and the corresponding  $(k, s)$ -tuple graph  $T_s^k(G)$ , we define the following variant of the 1-WL, taking into account edge labels, using the coloring

$$C_{i+1}^{1,\delta,*}(v_t) = (C_i^{1,\delta,*}(v_t), \{(C_i^{1,\delta,*}(v_s), \ell(v_t, v_s)) \mid v_s \in \delta(v_t)\}) \quad (12)$$

for  $i > 0$  and  $C_0^{1,\delta,*}(v_t) = \tau_{G[\mathbf{t}]}$  for  $v_t$  in  $V_T$ . Note that 1-WL and the variant defined via Equation (12) have the same asymptotic running time. The following lemma states that the  $(k, s)$ -LWL can be simulated on the  $(k, s)$ -tuple graph using the above variant of 1-WL.

**Lemma 6.** Let  $G$  be a graph,  $k \geq 1$ , and  $s$  in  $[k]$ , then

$$C_i^{k,s}(\mathbf{t}) = C_i^{k,s}(\mathbf{u}) \iff C_i^{1,\delta,*}(v_t) = C_i^{1,\delta,*}(v_u),$$

for all  $i \geq 0$ , and all  $(k, s)$ -tuples  $\mathbf{t}$  and  $\mathbf{u}$  in  $V(G)_s^k$ .

*Proof sketch.* Induction on the number of iterations using Definition 5.  $\square$

The *unrolling* of a neighborhood around a node of a given graph to a tree is defined as follows, see Figure 2 for an illustration.

**Definition 7.** Let  $G = (V, E, \ell)$  be a labeled (directed) graph and let  $v$  be in  $V$ . Then  $U_{G,v}^i = (W_i, F_i, l_i)$  for  $i \geq 0$  denotes the *unrolled tree*  $G$  around  $v$  at depth  $i$ , where

$$W_i = \begin{cases} \{v_{(0,v)}\} & \text{if } i = 0 \\ W_{i-1} \cup \{u_{(i,w_{(i-1,p)})} \mid u \in \delta(w) \text{ for } w_{(i-1,p)} \in W_{i-1}\} & \text{otherwise,} \end{cases}$$

and

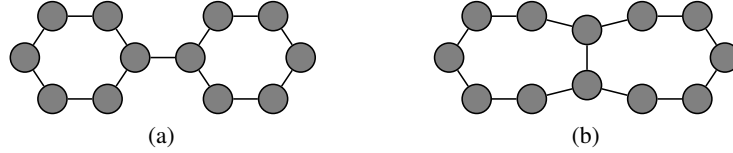
$$F_i = \begin{cases} \emptyset & \text{if } i = 0 \\ F_{i-1} \cup \{(w_{(i-1,p)}, u_{(i,w)}) \mid u \in \delta(w) \text{ for } w_{(i-1,p)} \in W_{i-1}\} & \text{otherwise.} \end{cases}$$

The label function is defined as  $l_i(u_{(j,p)}) = \ell(u)$  for  $u$  in  $V$ , and  $l_i(u_{(j,w)}) = \ell((w, u))$ . For notational convenience, we usually omit the subscript  $i$ .

In the following, we use the unrolled tree for the above defined  $(k, s)$ -tuple graph. For  $k \geq 2$  and  $s$  in  $[k]$ , we denote the *directed*, unrolled tree of the  $(k, s)$ -tuple graph of  $G$  around the node  $v_t$  at depth  $i$  for the tuple  $\mathbf{t}$  in  $V(G)_s^k$  by  $\mathbf{U}_{T_s^k(G), v_t}^i$ . For notational convenience, we write  $\mathbf{U}_{T, v_t}^i$ . Further, for two  $(k, s)$ -tuples  $\mathbf{t}$  and  $\mathbf{u}$ , we write

$$\mathbf{U}_{T, v_t}^i \simeq_{v_t \rightarrow v_u} \mathbf{U}_{T, v_u}^i \quad (13)$$



Figure 3: The graphs  $A_{k+2}$  and  $B_{k+2}$  for  $k = 4$ .

if there exists a (labeled) isomorphism  $\varphi$  between the two unrolled trees, mapping the (root) node  $v_t$  to  $v_u$ . Moreover, we need the following two results. The first one states that the 1-WL can distinguish any two directed, labeled non-isomorphic trees.

**Theorem 8** ((Busacker & Saaty, 1965; Valiente, 2002)). The 1-WL distinguishes any two directed, labeled non-isomorphic trees.

Using the first result, the second one states that the  $(k, s)$ -LWL can be simulated by the variant of the 1-WL of Equation (12) on the unrolled tree of the  $(k, s)$ -tuple graph, and hence can be reduced to a tree-isomorphism problem.

**Lemma 9.** Let  $G$  be a *connected* graph, then the  $(k, s)$ -LWL colors the tuples  $\mathbf{t}$  and  $\mathbf{u}$  in  $V(G)_s^k$  equally if and only if the corresponding unrolled  $(k, s)$ -tuple trees are isomorphic, i.e.,

$$C_i^{k,s}(\mathbf{t}) = C_i^{k,s}(\mathbf{u}) \iff \mathbf{U}_{T,v_t}^i \simeq_{v_t \rightarrow v_u} \mathbf{U}_{T,v_u}^i,$$

for all  $i \geq 0$ .

*Proof sketch.* First, by Lemma 6, we can simulate the  $(k, s)$ -LWL for the graph  $G$  using the  $k$ -tuple graph  $T^k(G)$ . Secondly, consider a node  $v_t$  in the  $k$ -tuple graph  $T^k(G)$  and a corresponding node in the unrolled tree around  $v_t$ . Observe that the neighborhoods for both nodes are identical. By definition, this holds for all nodes (excluding the leaves) in the unrolled tree. Hence, by Lemma 6, we can simulate the  $(k, s)$ -LWL for each tuple  $\mathbf{t}$  by running the 1-WL in the unrolled tree around  $v_t$  in the  $k$ -tuple graph. Since the 1-WL decides isomorphism for trees, see Theorem 8, the result follows.  $\square$

The following lemma shows that  $(k+1, 1)$ -LWL is strictly more expressive than  $(k, 1)$ -LWL for every  $k \geq 2$ .

**Lemma 10.** Let  $k \geq 2$ . Let  $G := C_{2(k+2)}$  and  $H := C_{(k+2)} \dot{\cup} C_{(k+2)}$ . Then, the graphs  $G$  and  $H$  are distinguished by  $(k+1, 1)$ -LWL, but they are not distinguished by  $(k, 1)$ -LWL.

*Proof.* We first show that the  $(k+1, 1)$ -LWL distinguishes the graphs  $G$  and  $H$ . Let  $\mathbf{v}$  in  $V(G)_1^{k+1}$  be a tuple  $(v_1, \dots, v_{k+1})$  such that  $v_1, \dots, v_{k+1}$  is a path of length  $k$  in  $G$ . Let  $\mathbf{w}$  in  $V(H)_1^{k+1}$  be a tuple  $(w_1, \dots, w_{k+1})$  such that  $w_1, \dots, w_{k+1}$  is a path of length  $k$  in  $H$ . By the structure of  $H$ , there exists a vertex  $w_{k+2}$  in  $V(H)$  such that  $w_1, \dots, w_{k+2}$  forms a cycle of length  $k+2$ . We claim that  $\mathbf{v}$  does not have any local 1-neighbor  $\mathbf{x}$  in  $V(G)_1^{k+1}$  such that  $\mathbf{x}$  is non-repeating, i.e., every vertex in  $\mathbf{x}$  is distinct. This holds because replacing the first vertex of  $\mathbf{v}$  by any other vertex of  $G$  will yield a disconnected tuple. On the other hand,  $\mathbf{w}$  admits a non-repeating, local 1-neighbor, obtained by replacing the first vertex  $w_1$  by  $w_{k+2}$ . Hence,  $(k+1, 1)$ -LWL distinguishes  $G$  and  $H$ .

Next, we show that  $(k, 1)$ -LWL cannot distinguish the graphs  $G$  and  $H$ . Indeed, for every  $j$  in  $[k]$ , every  $k$ -tuple  $\mathbf{x}$  in  $V(G)_1^k$  or  $V(H)_1^k$  has exactly two local  $j$ -neighbors, corresponding to the two neighbors  $y, z$  of the vertex  $x_j$ . The exact number of local  $j$ -neighbors of  $\mathbf{x}$  which additionally lie in  $V(G)_1^k$  (or  $V(H)_1^k$ ) depends purely only on the atomic type of  $\mathbf{x}$ , since the length of cycles in  $G$  and  $H$  is at least  $k+2$ . Hence, the  $(k, 1)$ -LWL neighborhood of every tuple in  $G$  or  $H$  depends only on its atomic type. This implies that  $(k, 1)$ -LWL does not refine the initial coloring for  $G$  as well as  $H$ , and hence it does not distinguish  $G$  and  $H$ .  $\square$

Although Lemma 10 already implies Theorem 4, the construction hinges on the fact that the graphs  $G$  and  $H$  are not connected. To address this, for  $k \geq 2$ , we introduce two connected graphs  $A_{k+2}$  and  $B_{k+2}$

defined as follows. The graph  $A_{k+2}$  has  $2(k+2)$  nodes and  $2(k+2)+1$  edges, and consists of two disjoint cycles on  $k+2$  nodes connected by a single edge. The graph  $B_{k+2}$  also has the same number of nodes and edges, and consists of two cycles on  $k+3$  nodes, each, sharing exactly two adjacent nodes. See Figure 3 for an illustration of the graphs  $A_{k+2}$  and  $B_{k+2}$  for  $k=4$ . We obtain the following result for the two graphs.

**Lemma 11.** Let  $k \geq 2$ , then the  $(k, 1)$ -LWL cannot distinguish the graphs  $A_{k+2}$  and  $B_{k+2}$ , while the  $(k+1, 1)$ -LWL can.

*Proof.* We first show the second part, i.e., the  $(k+1, 1)$ -LWL distinguishes the graphs  $A_{k+2}$  and  $B_{k+2}$ . Without loss of generality, assume that  $V(A_{k+2}) = \{a_1, \dots, a_{2(k+2)}\}$ , where  $(a_i, a_{i+1}) \in E(A_{k+2})$  for  $1 \leq i \leq k+1$ ,  $(a_1, a_{(k+2)}) \in E(A_{k+2})$ ,  $(a_i, a_{i+1}) \in E(A_{k+2})$  for  $(k+3) \leq i \leq 2(k+2)-1$  and  $(a_{k+3}, a_{2(k+2)}) \in E(A_{k+2})$ . The two cycles are connected by the edge  $(a_{(k+2)}, a_{(k+3)}) \in E(A_{k+2})$ . Further, assume  $V(B_{k+2}) = \{b_1, \dots, b_{2(k+2)}\}$  where  $(b_i, b_{i+1}) \in E(B_{k+2})$  for  $1 \leq i \leq 2(k+2)-1$  and  $(b_1, b_{2(k+2)+2}) \in E(B_{k+2})$ . Finally, the edge  $(b_1, b_{k+3}) \in E(B_{k+2})$  is shared by the two cycles of length  $k+3$  each.

Now, let  $\mathbf{t} = (a_1, \dots, a_{k+1}) \in V(A_{k+2})_1^{k+1}$ . Observe that the tuple  $(a_1, \dots, a_{k+2})$  is a  $(k+1)$ -neighbor of the tuple  $\mathbf{t}$ , inducing a graph on  $k+1$  nodes. Further, since the two cycles in the graph  $B_{k+2}$  have length  $k+3$ , there is no tuple without repeated nodes, in the graph that has a  $(k+1)$ -neighbor without repeated nodes. Hence, the two graphs are distinguished by  $(k+1, 1)$ -LWL.

We now show that the  $(k, 1)$ -LWL does not distinguish the graphs  $A_{k+2}$  and  $B_{k+2}$ . First, we construct a bijection  $\theta: V(A_{k+2}) \rightarrow V(B_{k+2})$  as induced by the following coloring:



Based on the bijection  $\theta$ , we define the bijection  $\theta^k: V(A_{k+2})_1^k \rightarrow V(B_{k+2})_1^k$ , by applying  $\theta$  component-wise to  $(k, s)$ -tuples. Observe that  $G[\mathbf{s}] \simeq G[\theta^k(\mathbf{s})]$  for  $\mathbf{s}$  in  $V(A_{k+2})_1^k$ .

**Claim 12.** Let  $\mathbf{s}$  be a tuple in  $V(G)_1^k$  and  $\mathbf{t} = \theta^k(\mathbf{s})$  in  $V(H)_1^k$ . Let  $N_j(\mathbf{s})$  and  $N_j(\mathbf{t})$  be the  $j$ -neighbors of the tuple  $\mathbf{s}$  and  $\mathbf{t}$ , respectively, for  $j$  in  $[k]$ . Then  $\theta^k$  yields a one-to-one correspondence between  $N_j(\mathbf{s})$  and  $N_j(\mathbf{t})$ . Consequently,  $G[\mathbf{u}] \simeq G[\theta^k(\mathbf{u})]$  for  $\mathbf{u}$  in  $N_j(\mathbf{s})$  and  $\theta^k(\mathbf{u})$  in  $N_j(\mathbf{t})$ .

*Proof.* For brevity, let  $F = A_{k+2}$  and  $K = B_{k+2}$ . The desired claim follows by observing that the bijective map  $\theta: V(A_{k+2}) \rightarrow V(B_{k+2})$  preserves neighborhoods, i.e. for every  $x$  in  $V(A_{k+2})$ ,  $\theta(N_F(x)) = N_K(\theta(x))$ .  $\square$

We now again leverage the above claim to show that  $C_i^{k,s}(\mathbf{s}) = C_i^{k,s}(\theta^k(\mathbf{s}))$  for  $i \geq 0$ , implying the required result. By a straightforward inductive argument, using Claim 12, we can inductively construct a tree isomorphism between the unrolled trees around the node  $v_s$  and  $v_t$  in the corresponding  $(k, s)$ -tuple graph such that  $\mathbf{U}_{T, v_s}^i \simeq_{v_s \rightarrow v_t} \mathbf{U}_{T, v_t}^i$ . By Lemma 9, this implies  $C_i^{k,s}(\mathbf{s}) = C_i^{k,s}(\theta(\mathbf{t}))$  for  $i \geq 0$ . This shows that  $(k, s)$ -LWL does not distinguish  $A_{k+2}$  and  $B_{k+2}$ . Hence, proved.  $\square$

Hence, the proof of Theorem 4 directly follows from the above Lemma 11. Moreover, we also show that the  $(k, 2)$ -LWL is more expressive than the  $(k, 1)$ -LWL.

**Proposition 13.** Let  $k \geq 2$ , then

$$(k, 2)\text{-LWL} \sqsubset (k, 1)\text{-LWL}.$$

*Proof.* As in Lemma 10, let  $G := C_{2(k+2)}$  and  $H := C_{(k+2)} \dot{\cup} C_{(k+2)}$ . By Lemma 10,  $G$  and  $H$  are not distinguished by the  $(k, 1)$ -LWL for  $k \geq 2$ . We claim that the  $(k, 2)$ -LWL distinguishes  $G$  and  $H$  for  $k = 2$  already. Since  $(k, 2)$ -LWL is at least as powerful as  $(2, 2)$ -LWL, this yields the desired claim.

With respect to  $(2, 2)$ -LWL, observe that the  $(2, 2)$ -tuple graph  $T_2^2(H)$  consists of four connected components while the  $(2, 2)$ -tuple graph  $T_2^2(G)$  consists of a single connected component. More precisely,

there exist two connected components of  $T_2^2(H)$  consisting only of 2-tuples containing two non-adjacent nodes, not being in the same connected component of the graph  $H$ . Note that none of these 2-tuples is adjacent to any 2-tuples in  $V(H)_1^2$ . Moreover, there exists no such connected component in  $T_2^2(G)$ . Moreover, note that the number of each neighbors of each 2-tuple of both graphs is exactly 4, excluding self loops. Hence, the  $(2, 2)$ -LWL will distinguish the two graphs.  $\square$

Moreover, the following result shows that increasing the parameter  $s$  results in higher expressivity. Formally, we show that the  $(k, k)$ -LWL is strictly more expressive than the  $(k, 2)$ -LWL. Note that we use vertex-colored graphs (rather than simple undirected graphs) in our proofs.

**Theorem 14.** Let  $k \geq 2$ , then

$$(k, k)\text{-LWL} \sqsubset (k, 2)\text{-LWL}.$$

For the proof of Theorem 14, we modify the construction employed in (Morris et al., 2020b), Appendix C.1.1., where they provide an infinite family of graphs  $(G_k, H_k)$ ,  $k$  in  $\mathbb{N}$ , such that (a)  $k$ -WL does not distinguish  $G_k$  and  $H_k$ , although (b)  $\delta$ - $k$ -LWL distinguishes  $G_k$  and  $H_k$ . Since our proof closely follows theirs, let us recall some relevant definitions from their paper.

*Construction of  $G_k$  and  $H_k$ .* Let  $K$  denote the complete graph on  $k + 1$  vertices (without any self-loops). The vertices of  $K$  are indexed from 0 to  $k$ . Let  $E(v)$  denote the set of edges incident to  $v$  in  $K$ : clearly,  $|E(v)| = k$  for all  $v$  in  $V(K)$ . We call the elements of  $V(K)$  as *base vertices*, and the elements of  $E(K)$  as *base edges*. Define the graph  $G_k$  as follows:

1. For the vertex set  $V(G_k)$ , we add
  - (a)  $(v, S)$  for each  $v$  in  $V(K)$  and for each *even* subset  $S$  of  $E(v)$ ,
  - (b) two vertices  $e^1, e^0$  for each edge  $e$  in  $E(K)$ .
2. For the edge set  $E(G_k)$ , we add
  - (a) an edge  $\{e^0, e^1\}$  for each  $e$  in  $E(K)$ ,
  - (b) an edge between  $(v, S)$  and  $e^1$  if  $v$  in  $e$  and  $e$  in  $S$ ,
  - (c) an edge between  $(v, S)$  and  $e^0$  if  $v$  in  $e$  and  $e$  not in  $S$ ,

For every  $v$  in  $K$ , the set of vertices of the form  $(v, S)$  is called the *vertex-cloud* for  $v$ . Similarly, for every edge  $e$  in  $E(K)$ , the set of vertices of the form  $\{e^0, e^1\}$  is called the *edge-cloud* for  $e$ .

Define a companion graph  $H_k$ , in a similar manner to  $G_k$ , with the following exception: in Step 1(a), for the vertex 0 in  $V(K)$ , we choose all *odd* subsets of  $E(0)$ . Counting vertices, we find that  $|V(G)| = |V(H)| = (k + 1) \cdot 2^{k-1} + \binom{k}{2} \cdot 2$ . This finishes the construction of graphs  $G$  and  $H$ . We set  $G_k := G$  and  $H_k := H$ .

*Distance two cliques.* A set  $S$  of vertices is said to form a *distance-two-clique* if the distance between any two vertices in  $S$  is exactly two. The following results were shown in (Morris et al., 2020b).

**Lemma 15** ((Morris et al., 2020b)). The following holds for graphs  $G_k$  and  $H_k$  defined above.

- There exists a distance-two-clique of size  $(k + 1)$  inside  $G_k$ .
- There does not exist a distance-two-clique of size  $(k + 1)$  inside  $H_k$ .

Hence,  $G_k$  and  $H_k$  are non-isomorphic.

**Lemma 16** ((Morris et al., 2020b)). The  $\delta$ - $k$ -LWL distinguishes  $G_k$  and  $H_k$ . On the other hand,  $k$ -WL does not distinguish  $G_k$  and  $H_k$ .

We are ready to present the proof of Theorem 14. The crux is to subdivide the edges in  $G_k$  and  $H_k$  into sufficiently long paths. The local nature of the  $(k, 2)$ -LWL then ensures that a tuple considered by  $(k, 2)$ -LWL actually accesses only a constant number, in fact, just two, of original vertices of  $G_k$  and  $H_k$ . As a result, it is not possible for the  $(k, 2)$ -LWL to aggregate the right “ $k$ -ary” information necessary for distinguishing such graphs. We proceed with the details.

*Proof of Theorem 14.* Observe that the  $(k, k)$ -LWL is the same as the  $\delta$ - $k$ -LWL. Hence, it suffices to show an infinite family of graphs  $(X_k, Y_k)$ ,  $k$  in  $\mathbb{N}$ , such that (a)  $(k, 2)$ -LWL does not distinguish  $X_k$  and  $Y_k$ , although (b)  $\delta$ - $k$ -LWL distinguishes  $X_k$  and  $Y_k$ .

Let  $X_k$  be the graph obtained from the graph  $G_k$  as follows. First, for every base vertex  $v$  in  $V(K)$ , every vertex of  $V(G_k)$  in the vertex cloud for  $v$  receives a color  $\text{Red}_v$ . Hence, vertex-clouds form color classes, where each such class has a distinct color. Similarly, for every base edge  $e$  in  $E(K)$ , every vertex of  $V(G_k)$  in the edge cloud for  $e$  receives a color  $\text{Blue}_e$ . Finally, let  $\Delta > 3k$ . Then, we replace every edge  $e$  in  $G_k$  by a path of length  $\Delta$ , such that every vertex on this path is colored with the color  $(\{c, c'\})$ , where  $c$  and  $c'$  are the colors of the endpoints of  $e$  in  $G_k$ . We call such path vertices *auxiliary* vertices. The graph  $Y_k$  is obtained from  $H_k$  by an identical construction.

First, we show that the  $(k, 2)$ -LWL does not distinguish graphs  $X_k$  and  $Y_k$ . In fact, we show a stronger statement: the  $(k, 2)$ -LWL does not even refine the initial coloring of graphs  $X_k$  and  $Y_k$ . Indeed, let  $\mathbf{x}$  be a  $k$ -tuple in  $V(X_k)_2^k$ . Let  $\text{set}(\mathbf{x})$  be the set of vertices present in  $\mathbf{x}$ . Observe that (a)  $\mathbf{x}$  induces at most two connected components, and (b) the cardinality of  $\text{set}(\mathbf{x}) \leq k$ . There are three further possibilities, as follows.

1. The tuple  $\mathbf{x}$  induces exactly one component. Since  $\Delta > 3k$ , there is at most one vertex-cloud or edge-cloud which intersects with  $\text{set}(\mathbf{x})$ . Moreover, exactly one vertex in such a cloud belongs to  $\text{set}(\mathbf{x})$ . This holds because any two vertices in a cloud are at a distance greater than  $\Delta$  apart. In this case, the local neighborhood of  $\mathbf{x}$  used by  $(k, 2)$ -LWL depends solely on the atomic type of  $\mathbf{x}$ .
2. The tuple  $\mathbf{x}$  induces two components  $C_1$  and  $C_2$  which are at distance at least 3 from each other. Since  $\Delta > 3k$ , there is at most one vertex-cloud or edge-cloud which intersects with either component. Moreover, exactly one vertex in such a cloud belongs to  $\text{set}(\mathbf{x})$ . This holds because any two vertices in a cloud are at a distance  $> \Delta$  apart. In this case as well, the local neighbourhood of  $\mathbf{x}$  used by  $(k, 2)$ -LWL depends solely on the atomic type of  $\mathbf{x}$ .
3. The tuple  $\mathbf{x}$  induces two components  $C_1$  and  $C_2$  which are at distance at most 2 from each other. Since  $\Delta > 3k$ , one of these two components must consist entirely of auxiliary vertices. Hence, again, the local neighbourhood of  $\mathbf{x}$  used by  $(k, 2)$ -LWL depends solely on the atomic type of  $\mathbf{x}$ .

Since in all these cases, the local neighborhood of a tuple depends only on its atomic type, it suffices to verify that  $V(X_k)_2^k$  and  $V(Y_k)_2^k$  have the same histogram of atomic types. This holds for the following reason. The histogram of atomic types in the padded graphs  $X_k$  and  $Y_k$  is equal if and only if the histogram of atomic types in the original graphs  $G_k$  and  $H_k$  is equal. If  $V(X_k)_2^k$  and  $V(Y_k)_2^k$  have a different histogram of atomic types, then  $V(G_k)_2^k$  and  $V(H_k)_2^k$  also have a different histogram of atomic types. This implies that  $k$ -WL distinguishes  $G_k$  and  $H_k$  in the very first round, which contradicts Lemma 16. Hence,  $(k, 2)$ -LWL does not progress beyond the initial coloring and fails to distinguish  $X_k$  and  $Y_k$ .

Next, we show that the  $\delta$ - $k$ -LWL distinguishes graphs  $X_k$  and  $Y_k$ . Our proof closely follows the corresponding proof in (Morris et al., 2020b). Instead of showing a discrepancy in the number of distance-two-cliques, we instead use colored-distance- $(2\Delta + 1)$ -cliques defined as follows. Let  $S$  be a set of vertices belonging to the vertex clouds. The set  $S$  is said to form a colored-distance- $(2\Delta + 1)$ -clique if any two vertices in  $S$  are connected by a path of exactly  $2\Delta + 1$  vertices, of which  $2\Delta$  are auxiliary vertices and one vertex is a vertex from an edge-cloud. Analogous to their proof, it can be shown that (a) there exists a colored-distance- $(2\Delta + 1)$ -clique of size  $(k + 1)$  inside  $X_k$ , and (b) there does not exist a colored-distance- $(2\Delta + 1)$ -clique of size  $(k + 1)$  inside  $Y_k$ , and hence, (c)  $X_k$  and  $Y_k$  are non-isomorphic. Finally, we claim that the  $\delta$ - $k$ -LWL is powerful enough to detect colored-distance- $(2\Delta + 1)$ -cliques. The proof is analogous to that of (Morris et al., 2020b). This yields that  $\delta$ - $k$ -LWL distinguishes graphs  $X_k$  and  $Y_k$ .  $\square$

## C.2 ASYMPTOTIC RUNNING TIME

In the following, we bound the asymptotic running time of the  $(k, s)$ -LWL. Due to Lemma 6, we can upper-bound the running time of  $(k, s)$ -LWL for a given graph by upper-bounding the time to construct the  $(k, s)$ -tuple graph and running the 1-WL variant of Equation (12) on top. Proposition 18 establishes an upper bound on the asymptotic running time for constructing the  $(k, s)$ -tuple graph from a given graph. Thereto, we assume a  $d$ -bounded degree graph  $G$ , for  $d \geq 1$ , i.e., each node has at most  $d$  neighbors.

To prove the proposition, we define  $(k, s)$ -multisets. Let  $G$  be a graph,  $k \geq 1$ , and  $s \in [k]$ , then the set of  $(k, s)$ -multisets

$$S(G)_s^k = \{\{v_1, \dots, v_k\} \mid \mathbf{v} \in V(G)_s^k\}$$

contains the set of multisets inducing subgraphs of  $G$  on at most  $k$  nodes with at most  $s$  components. The following results upper-bounds the running time for the construction of  $S(G)_s^k$ .

**Proposition 17.** Let  $G$  be a  $d$ -bounded degree graph,  $k \geq 2$ , and  $s$  in  $[k - 1]$ . Then Algorithm 1 computes  $S_s^k(G)$  in  $\tilde{O}(n^s \cdot k^{k-s} (d+1)^{k-s})$

*Proof.* Let  $T'$  be an element in  $S(G)_s^{c+1}$  for  $c+1 \leq k$ . By definition of  $S(G)_s^c$  and  $S(G)_s^{c+1}$ , there exists a  $c$ -element multiset  $T$  in  $S(G)_s^c$  such that  $T' = T \cup \{v\}$  is in  $S(G)_s^{c+1}$  for a node  $v$  in  $V(G)$ . Since  $s$  is fixed,  $v$  is either in the neighborhood  $\delta(w)$  for  $w$  in  $T$  or  $v = w'$  for  $w'$  in  $T$ . Hence, lines 7 to 9 in Algorithm 1 generate  $S(G)_s^{c+1}$  from  $S(G)_s^c$ . The set data structure  $R$  makes sure that the final solution will not contain duplicates. The running time follows directly when using, e.g., a red-black tree, to represent the set  $R$ .  $\square$

Based on the above result, we can easily construct  $T_s^k(G)$  from  $S_s^k(G)$ , implying the following result.

**Proposition 18.** Let  $G$  be a  $d$ -bounded degree graph,  $k \geq 3$ , and  $s \in [k - 1]$ . Then we can compute  $T_s^k(G)$  in  $\tilde{O}(n^s \cdot k^{k-s} (d+1)^{k-s+1} \cdot k! \cdot k)$

*Proof.* The running time follows directly from Proposition 17. That is, from  $S_s^k(G)$  we can generate the set  $V_s^k(G)$  by generating all permutations of each element in the former. By iterating over each resulting  $(k, s)$ -tuple and each component of such  $(k, s)$ -tuple, we can construct the needed adjacency information.  $\square$

Hence, unlike for  $k$ -WL, the running time of  $(k, s)$ -LWL does not depend on  $n^k$  for an  $n$ -node graph and is solely dictated by  $s, k$ , and the sparsity of the graph.

---

#### Algorithm 1 Generate $(k, s)$ -multisets

---

**Input:** Graph  $G, k, s$ , and  $S(G)_s^s$   
**Output:**  $(k, s)$ -multiset  $S(G)_s^k$

- 1: Let  $R$  be a empty set data structure
- 2: **for**  $M \in S(G)_s^s$  **do**
- 3:     Let  $S$  be a queue data structure containing only  $(M, s)$
- 4:     **while**  $S$  not empty **do**
- 5:         Pop  $(T, c)$  from queue  $S$
- 6:         **if**  $c+1 \leq k$  **then**
- 7:             **for**  $t \in T$  **do**
- 8:                 **for**  $u \in \delta(t) \cup \{t\}$  **do**
- 9:                     Add  $(T \cup \{u\}, c+1)$  to  $S$
- 10:         **else**
- 11:             Add  $T$  to  $R$
- 12: **return**  $R$

---

Moreover, observe that the upper bound given in Proposition 18, by leveraging Lemma 6, also upper-bounds the asymptotic running time for one iteration of the  $(k, s)$ -LWL.

## D SPEQNETS (EXTENDED)

The following result demonstrates the expressive power of  $(k, s)$ -SpeqNets, in terms of distinguishing non-isomorphic graphs.

**Theorem 19.** Let  $(V, E, \ell)$  be a labeled graph, and let  $k \geq 1$  and  $s$  in  $[k]$ . Then for all  $t \geq 0$ , there exists a sequence of weights  $\mathbf{W}^{(t)}$  such that

$$C_t^{k,s}(\mathbf{v}) = C_t^{k,s}(\mathbf{w}) \iff f^{(t)}(\mathbf{v}) = f^{(t)}(\mathbf{w}).$$

Hence, the following holds for all  $k \geq 1$ :

$$(k, s)\text{-SpeqNet} \equiv (k, s)\text{-LWL}.$$

*Proof sketch.* First, observe that the  $(k, s)$ -LWL can be simulated on an appropriate node- and edge-labeled graph, see Lemma 6. Secondly, following the proof of (Morris et al., 2019, Theorem 2), there exists a parameter matrix  $W_2^{(t)}$  such that we can injectively map each multiset in Equation (3), representing the local  $j$ -neighbors for  $j$  in  $[k]$ , to a  $d$ -dimensional vector. Moreover, we concatenate  $j$  to each such vector to distinguish between different neighborhoods. Again, by (Morris et al., 2019, Theorem 2), there exists a parameter matrix  $W_1^{(t)}$  such that we can injectively map the set of resulting  $k$  vectors to a unique vector representation. Alternatively, one can concatenate the resulting  $k$  vectors and use a multi-layer perceptron to learn a joint, lower-dimensional representation.  $\square$

Note that it is not possible to come up with an architecture and weight assignments of  $f_{\text{mg}}^{W_1}$  and  $f_{\text{agg}}^{W_2}$ , such that it becomes more powerful than the  $(k, s)$ -LWL, see (Morris et al., 2019). However, all results from the previous section can be lifted to the neural setting. Analogously to GNNs, the above architecture can naturally handle continuous node and edge labels. By using the tools developed in (Azizian & Lelarge, 2020), it is straightforward to show that the above architecture is universal, i.e., it can approximate any continuous, bounded, permutation-invariant function over graphs up to an arbitrarily small additive error.

#### D.1 NODE-, EDGE-, AND SUBGRAPH-LEVEL LEARNING TASKS

The above architecture computes representations for  $(k, s)$ -tuples, making it mostly suitable for graph-level learning tasks, e.g., graph classification or regression. However, it is also possible to derive neural architectures based on the  $(k, s)$ -LWL for node- and edge-level learning tasks, e.g., node or link prediction. Given a graph  $G$ , to learn a node feature for node  $v$ , we can simply pool over the features learned for  $(k, s)$ -tuples containing the node  $v$  as a component. That is, let  $t > 0$ , then we consider the multisets

$$m^t(v)_i = \{ \{ f^{(t-1)}(\mathbf{t}) \mid \mathbf{t} \in V(G)_s^k \text{ and } t_i = v \} \} \quad (14)$$

for  $i$  in  $[k]$ . Hence, to compute a vectorial representation of the node  $v$ , we compute a vectorial representation of  $m^t(v)_i$  for  $i$  in  $[k]$ , e.g., using a neural architecture for multi-sets, see (Wagstaff et al., 2021), followed by learning a joint vectorial representation for the node  $v$ . Again, by (Azizian & Lelarge, 2020), it is straightforward to show that the above architecture is universal, i.e., it can approximate any continuous, bounded, permutation-equivariant function over graphs up to an arbitrarily small additive error. Note that the above approach can be directly generalized to learn subgraph representations on an arbitrary number of vertices.

## E DETAILS ON EXPERIMENTS

**Datasets** To compare the  $(k, s)$ -LWL-based kernels, we used the well-known graph classification benchmark datasets from (Morris et al., 2020a), see Table 3 for dataset statistics and properties.<sup>5</sup> To compare the  $(k, s)$ -SpeqNet architecture to GNN baselines, we used the ALCHEMY (Chen et al., 2019a) and the QM9 (Ramakrishnan et al., 2014; Wu et al., 2018) graph regression datasets, again see Table 1 for dataset statistics and properties. Following Morris et al. (2020b), we opted for not using the 3D-coordinates of the ALCHEMY dataset to solely show the benefits of the (sparse) higher-order structures concerning graph structure and discrete labels. To investigate the performance of the architecture for node classification, we used the WEBKB datasets (Pei et al., 2020), see Table 4 for dataset statistics and properties.

**Kernels** We implemented the  $(k, s)$ -LWL and  $(k, s)$ -LWL<sup>+</sup> for  $k$  in  $\{2, 3\}$  and  $s$  in  $\{1, 2\}$ . We compared our kernels to the Weisfeiler–Leman subtree kernel (1-WL) (Shervashidze et al., 2011), the Weisfeiler–Leman Optimal Assignment kernel (WLOA) (Kriege et al., 2016), the graphlet kernel (GR) (Shervashidze et al., 2009), and the shortest-path kernel (Borgwardt & Kriegel, 2005) (SP). Further, we implemented the higher-order kernels  $\delta$ - $k$ -LWL,  $\delta$ - $k$ -LWL<sup>+</sup>,  $\delta$ - $k$ -WL, and  $k$ -WL kernel for  $k$  in  $\{2, 3\}$  as outlined in (Morris et al., 2020b). All kernels were (re-)implemented in C++11. For the graphlet kernel, we counted (labeled) connected subgraphs of size 3. We followed the evaluation guidelines outlined in (Morris et al., 2020a).

<sup>5</sup>All datasets are publicly available at [www.graphlearning.io](http://www.graphlearning.io).

**Neural architectures** We used the GIN and GIN- $\epsilon$  architecture (Xu et al., 2019) as neural baselines. For data with (continuous) edge features, we used a 2-layer MLP to map them to the same number of components as the node features and combined them using summation (GINE and GINE- $\epsilon$ ). For the evaluation of the  $(k, s)$ -SpeqNet neural architectures, we implemented them using PYTORCH GEOMETRIC (Fey & Lenssen, 2019), using a Python-wrapped C++11 preprocessing routine to compute the computational graphs for the higher-order GNNs. We used the GIN- $\epsilon$  layer to express  $f_{\text{mrg}}^{W_1}$  and  $f_{\text{agg}}^{W_2}$  of Equation (3). For the GNN baseline for the QM9 dataset, following (Gilmer et al., 2017), we used a complete graph, computed pairwise  $\ell_2$  distances based on the 3D coordinates, and concatenated them to the edge features. We note here that our intent is not the best state-of-the-art, physical knowledge-incorporating architectures, e.g., DimeNet (Klicpera et al., 2020) or Cormorant (Anderson et al., 2019), but to solely show the benefits of the local, sparse higher-order architectures compared to the corresponding (1-dimensional) GNN. For the  $(k, s)$ -SpeqNet architectures, in the case of the QM9 dataset, to compute the initial features, for each  $(k, s)$ -tuple, we concatenated the node and edge features, computed pairwise  $\ell_2$  distances based on the 3D coordinates, and a one-hot encoding of the (labeled) isomorphism type. Finally, we used a 2-layer MLP to learn a joint, initial vectorial representation. For the node classification experiments, we used mean pooling to implement Equation (14) and a standard GCN layer for all experiments, including the  $(k, s)$ -SpeqNet architectures. Further, we used the architectures (SDRF) outlined in (Topping et al., 2021) as baselines.

For the kernel experiments, we computed the (cosine) normalized Gram matrix for each kernel. We computed the classification accuracies using the  $C$ -SVM implementation of LIBSVM (Chang & Lin, 2011), using 10-fold cross-validation. We repeated each 10-fold cross-validation ten times with different random folds and report average accuracies and standard deviations.

Following the evaluation method proposed in (Morris et al., 2020a), the  $C$ -parameter was selected from  $\{10^{-3}, 10^{-2}, \dots, 10^2, 10^3\}$  using a validation set sampled uniformly at random from the training fold (using 10% of the training fold). Similarly, the numbers of iterations of the  $(k, s)$ -LWL,  $(k, s)$ -LWL<sup>+</sup>, 1-WL, WLOA,  $\delta$ - $k$ -LWL,  $\delta$ - $k$ -LWL<sup>+</sup>, and  $k$ -WL were selected from  $\{0, \dots, 5\}$  using the validation set. Moreover, for the  $(k, s)$ -LWL<sup>+</sup> and  $\delta$ - $k$ -LWL<sup>+</sup>, we only added the additional label function  $\#$  on the last iteration to prevent overfitting. We report computation times for the  $(k, s)$ -LWL,  $(k, s)$ -LWL<sup>+</sup>, WLOA,  $\delta$ - $k$ -LWL,  $\delta$ - $k$ -LWL<sup>+</sup>, and  $k$ -WL with five refinement steps. All kernel experiments were conducted on a workstation with 64GB of RAM using a single core. Moreover, we used the GNU C++ Compiler 7.4.0 with the flag `-O2`.

For comparing the kernel approaches to GNN baselines, we used 10-fold cross-validation and again used the approach outlined in (Morris et al., 2020a). The number of components of the (hidden) node features in  $\{32, 64, 128\}$  and the number of layers in  $\{1, 2, 3, 4, 5\}$  of the GIN and GIN- $\epsilon$  layer were again selected using a validation set sampled uniformly at random from the training fold (using 10% of the training fold). We used mean pooling to pool the learned node embeddings to a graph embedding and used a 2-layer MLP for the final classification, using a dropout layer with  $p = 0.5$  after the first layer of the MLP. We repeated each 10-fold cross-validation ten times with different random folds and report the average accuracies and standard deviations. Due to the different training methods, we do not provide computation times for the GNN baselines.

For the larger molecular regression tasks, ALCHEMY and QM9, we closely followed the hyperparameters found in (Chen et al., 2019a) and (Gilmer et al., 2017), respectively, for the GINE- $\epsilon$  layers. That is, we used six layers with 64 (hidden) node features and a set2seq layer (Vinyals et al., 2016) for graph-level pooling, followed by a 2-layer MLP for the joint regression of the twelve targets. We used the same architecture details and hyperparameters for the  $(k, s)$ -SpeqNet. For the ALCHEMY, we used the subset of 12 000 graphs as in (Morris et al., 2020b). For both datasets, we uniformly and at random sampled 80% of the graphs for training, and 10% for validation and testing, respectively. Moreover, following (Chen et al., 2019a; Gilmer et al., 2017), we normalized the targets of the training split to zero mean and unit variance. We used a single model to predict all targets. Following (Klicpera et al., 2020, Appendix C), we report mean standardized MAE and mean standardized logMAE. We repeated each experiment five times and report average scores and standard deviations. We used the provided ten train, validation, and test splits for the node classification datasets. All neural experiments were conducted on a workstation with one GPU card with 32GB of GPU memory.

To compare training and testing times between the  $(2, 1)$ -SpeqNet,  $(2, 2)$ -SpeqNet, GINE- $\epsilon$  architectures, we trained all three models on ALCHEMY (10K) and QM9 to convergence, divided by the number of epochs, and calculated the ratio with respect to the average epoch computation time of the  $(2, 1)$ -SpeqNet

(average computation time of dense or baseline layer divided by average computation time of the (2, 1)-SpeqNet).

## F ADDITIONAL EXPERIMENTAL RESULTS

**A3 Kernels** See Table 5. Clearly, for the same  $k$  and  $s < k$ , the  $(k, s)$ -LWL improves over the  $k$ -WL and its (local) variants. For example, on the ENZYMES dataset, the (2, 1)-LWL is more than 20 times faster in terms of computation times compared to the  $\delta$ -2-LWL. The speed up is even more significant for the non-local 2-WL. This speed-up factor increases more as  $k$  increases, e.g., the (3, 1)-LWL is more than 1 700 times faster compared to the 3-WL, whereas the (3, 2)-LWL is still more than 87 times faster, while giving better accuracies. Similar speed-up factors can be observed over all datasets.

*Neural architectures* See Table 6. The (2, 1)-SpeqNet severely speeds up the computation time across both datasets. Specifically, on the ALCHEMY dataset, the (2, 1)-SpeqNet is 1.3 times faster compared to the (2, 2)-SpeqNet, while requiring twice the computation time of the GINE- $\varepsilon$  but achieving a lower MAE. More interestingly, on the QM9 dataset, the (2, 1)-SpeqNet is 3.4 times faster compared to the (2, 2)-SpeqNet, while also being 1.3 times faster compared to the GINE- $\varepsilon$ . The speed-up over GINE- $\varepsilon$  is most likely due to the latter considering the complete graph to compute all pairwise  $\ell_2$  distances, whereas the (2, 1)-SpeqNet only considers connected node pairs.

Table 3: Dataset statistics and properties for graph-level prediction tasks,  $\dagger$ —Continuous vertex labels following Gilmer et al. (2017), the last three components encode 3D coordinates.

Dataset	Properties					
	Number of graphs	Number of classes/targets	$\varnothing$ Number of nodes	$\varnothing$ Number of edges	Node labels	Edge labels
ENZYMES	600	6	32.6	62.1	✓	✗
IMDB-BINARY	1 000	2	19.8	96.5	✗	✗
IMDB-MULTI	1 500	3	13.0	65.9	✗	✗
MUTAG	188	2	17.9	19.8	✓	✗
NCI1	4 110	2	29.9	32.3	✓	✗
PTC_FM	349	2	14.1	14.5	✓	✗
PROTEINS	1 113	2	39.1	72.8	✓	✗
REDDIT-BINARY	2 000	2	429.6	497.8	✗	✗
ALCHEMY	202 579	12	10.1	10.4	✓	✓
QM9	129 433	12	18.0	18.6	✓(13+3D) $\dagger$	✓(4)

Table 4: Dataset statistics and properties for node-level prediction tasks.

Dataset	Properties		
	Number of nodes	Number of edges	Number of node features
CORNELL	183	295	1 703
TEXAS	183	309	1 703
WISCONSIN	251	490	1 703



Table 5: Overall computation times for selected datasets in seconds (Number of iterations for WL-based methods: 5), OOT—Computation did not finish within one day (24h), OOM—Out of memory.

Graph Kernel		Dataset			
		ENZYMES	IMDB-BINARY	NCII	PROTEINS
Global	1-WL	<1	<1	2	<1
	2-WL	183	71	893	14 755
	3-WL	74 712	18 180	OOT	OOM
	$\delta$ -2-WL	294	89	1 469	14 620
	$\delta$ -3-WL	64 486	17 464	OOT	OOM
	Local	$\delta$ -2-LWL	20	22	92
$\delta$ -2-LWL <sup>+</sup>		22	23	103	177
$\delta$ -3-LWL		4 453	3 496	18 035	17 848
$\delta$ -3-LWL <sup>+</sup>		4 973	3 748	20 644	OOM
$(k, s)$ -LWL	(2, 1)-LWL	2	11	7	4
	(2, 1)-LWL <sup>+</sup>	2	12	7	5
	(3, 1)-LWL	36	871	72	87
	(3, 1)-LWL <sup>+</sup>	39	1 064	82	100
	(3, 2)-LWL	740	2 153	1 928	5 128
	(3, 2)-LWL <sup>+</sup>	1 097	2 797	2 837	6 754

Table 6: Average speed-up ratios over all epochs (training and testing).

Method	Dataset	
	ALCHEMY (10K)	QM9
GINE- $\epsilon$	0.5	1.3
(2, 1)-SpeqNet	1.0	1.0
(2, 2)-SpeqNet	1.3	3.4