
Gradual Weisfeiler-Leman: Slow and Steady Wins the Race

Anonymous Author(s)

Anonymous Affiliation

Anonymous Email

Abstract

The classical Weisfeiler-Leman algorithm aka color refinement is fundamental for graph learning and central for successful graph kernels and graph neural networks. Originally developed for graph isomorphism testing, the algorithm iteratively refines vertex colors. On many datasets, the stable coloring is reached after a few iterations and the optimal number of iterations for machine learning tasks is typically even lower. This suggests that the colors diverge too fast, defining a similarity that is too coarse. We generalize the concept of color refinement and propose a framework for gradual neighborhood refinement, which allows a slower convergence to the stable coloring and thus provides a more fine-grained refinement hierarchy and vertex similarity. We assign new colors by clustering vertex neighborhoods, replacing the original injective color assignment function. Our approach is used to derive new variants of existing graph kernels and to approximate the graph edit distance via optimal assignments regarding vertex similarity. We show that in both tasks, our method outperforms the original color refinement with only moderate increase in running time advancing the state of the art.

1 Introduction

The (1-dimensional) Weisfeiler-Leman algorithm, also referred to as *color refinement*, iteratively refines vertex colors by encoding colors of neighbors and was originally developed as a heuristic for the graph isomorphism problem. Although it cannot distinguish some non-isomorphic graph pairs, for example strongly regular graphs, it succeeds in many cases. It is widely used as a sub-routine in isomorphism algorithms today to reduce ambiguities that have to be resolved by backtracking search [1]. It has also gained high popularity in graph learning, where the technique is used to define graph kernels [2–5] and to formalize the expressivity of graph neural networks, see the recent surveys [6, 7]. Graph kernels based on Weisfeiler-Leman refinement provide remarkable predictive performance while being computationally highly efficient. The original Weisfeiler-Leman subtree kernel [2] and its variants and extensions, e.g., [3–5], provide state-of-the-art classification accuracy on many datasets and are widely used baselines. The update scheme of the Weisfeiler-Leman algorithm is similar to the idea of neighborhood aggregation in graph neural networks (GNNs). It has been shown that (i) the expressive power of GNNs is limited by the Weisfeiler-Leman algorithm, and (ii) that GNN architectures exist that reach this expressive power [8, 9].

As a consequence of its original application, the Weisfeiler-Leman algorithm assigns discrete colors and does not allow distinguishing minor or major differences in vertex neighborhood but considers two colors as either the same or different. Most Weisfeiler-Leman graph kernels match vertex colors of the first few refinement steps by equality, which can be considered as too rigid, since these colors encode complex neighborhood structures. In machine learning tasks, a more fine-grained differentiation appears promising. Often data is noisy, which in graphs can show for example in small differences in vertex degree. Such differences get picked up by the refinement strategy of the Weisfeiler-Leman algorithm and cannot be distinguished from significant differences.

We address this problem by providing a different approach to the refinement step of the Weisfeiler-Leman algorithm: We replace the injective relabeling function with a non-injective one, to gain a more gradual refinement of colors. This allows to obtain a finer vertex similarity measure, that can

44 distinguish between large and small changes in vertex neighborhoods with increasing radius. We
 45 characterize the set of functions that, while not necessarily injective, guarantee that the stable coloring
 46 of the original Weisfeiler-Leman algorithm is reached after a possibly higher number of iterations.
 47 Thus, our approach preserves the expressive power of the Weisfeiler-Leman algorithm. We discuss a
 48 possible realization of such a function and use k -means clustering in our experimental evaluation as
 49 an exemplary one.

50 **Our Contribution.**

- 51 1. We propose refining, neighborhood preserving (*renep*) functions, which generalize the concept
 52 of color refinement. This family of functions leads to the coarsest stable coloring while only
 53 incorporating direct neighborhoods.
- 54 2. We show the connections of our approach to the original Weisfeiler-Leman algorithm, as well as
 55 other vertex refinement strategies.
- 56 3. We propose two new graph kernels based on *renep* functions, that outperform state-of-the-art
 57 kernels on synthetic and real-world datasets, with only moderate increase in running time.
- 58 4. We apply our new approach for approximating the graph edit distance via bipartite graph
 59 matching and show that it outperforms state-of-the-art heuristics.

60 **2 Related Work**

61 Various graph kernels based on the standard Weisfeiler-Leman refinement have been proposed [2–5].
 62 Recent comprehensive experimental evaluations confirm their high classification accuracy on many
 63 real-world datasets [10, 11]. These approaches implicitly match colors by equality, which can be
 64 considered as too rigid, since colors encode unfolding trees representing complex neighborhood
 65 structures. Some recent works address this problem: Yanardag and Vishwanathan [12] introduced
 66 similarities between colors using techniques inspired by natural language processing, that were
 67 subsequently refined by Narayanan et al. [13]. Schulz et al. [14] define a distance function between
 68 colors by comparing the associated unfolding trees using a tree edit distance. Based on this distance
 69 the colors are clustered to obtain a new graph kernel. Although the tree edit distance is polynomial-
 70 time computable, the running time of the algorithm is very high. A kernel based on the Wasserstein
 71 distance of sets of unfolding trees was proposed by Fang et al. [15]. The vertices of the graphs are
 72 embedded into ℓ_1 space using an approximation of the tree edit distance between their unfolding
 73 trees. A graph can then be seen as a distribution over those embeddings. While the function proposed
 74 is not guaranteed to be positive semi-definite, the method showed results similar to and in some cases
 75 exceeding state-of-the-art techniques. The running time, however, is still very high and the method is
 76 only feasible for unfolding trees of small height.

77 These approaches define similarities between Weisfeiler-Leman colors and the associated unfolding
 78 trees. Our approach, in contrast, alters the Weisfeiler-Leman refinement procedure itself and does not
 79 rely on computationally expensive matching of unfolding trees.

80 **3 Preliminaries**

81 In this section we provide the definitions necessary to understand our new vertex refinement algorithm.
 82 We first give a short introduction to graphs and the original Weisfeiler-Leman algorithm, before we
 83 cover graph kernels.

84 **Graph Theory.** A graph $G = (V, E, \mu, \nu)$ consists of a set of vertices V , denoted by $V(G)$, a set
 85 of edges $E(G) = E \subseteq \binom{V}{2}$ between the vertices, a labeling function for the vertices $\mu: V \rightarrow L$,
 86 and a labeling function for the edges $\nu: E \rightarrow L$. The set L contains categorical labels, which can
 87 be represented as natural numbers. We discuss only undirected graphs and denote an edge between
 88 u and v by uv . The set of neighbors of a vertex $v \in V$ is denoted by $N(v) = \{u \mid uv \in E\}$. A
 89 (*rooted*) tree T is a simple (no self-loops or multi-edges), connected graph without cycles and with a
 90 designated root node r . A tree T' is a *subtree* of a tree T , denoted by $T' \subseteq T$, iff $V(T') \subseteq V(T)$.
 91 The root of T' is the node closest to the root in T .

92 A vertex coloring $c: V(G) \rightarrow \mathbb{N}_0$ of a graph G is a function assigning each vertex a color. For
 93 vertices with $c(u) = c(v)$ we also write $u \approx_c v$. A coloring π on a set S is a *refinement* of (or *refines*)

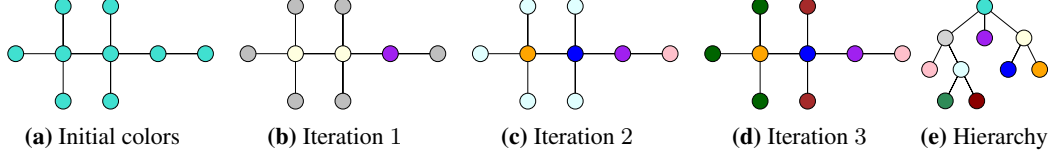


Figure 1: Initial coloring and results of the first three iterations of the Weisfeiler-Leman algorithm. To use less colors for this example, vertices with a unique color do not get a new color. The color hierarchy shows the development of the colors over the refinement iterations.

94 a coloring π' , iff $s_1 \approx_\pi s_2 \Rightarrow s_1 \approx_{\pi'} s_2$ for all s_1, s_2 in S . We denote this by $\pi \preceq \pi'$ and write
 95 $\pi \equiv \pi'$ if $\pi \preceq \pi'$ and $\pi' \preceq \pi$. If $\pi \preceq \pi'$ and $\pi \not\equiv \pi'$, we say that π is a *strict refinement* of π' , written
 96 $\pi \prec \pi'$. The refinement relation defines a partial ordering on the colorings.

97 **Color Hierarchy.** We consider a sequence of vertex colorings $(\pi_0, \pi_1, \dots, \pi_h)$ with $\pi_h \preceq \dots \preceq \pi_0$
 98 and assume that for $i \neq j$ the colors assigned by π_i and π_j are distinct. We can interpret such a
 99 sequence of colorings as a *color hierarchy*, i.e., a tree \mathcal{T}_h that contains a node for each color $c \in$
 100 $\{\pi_i(v) \mid i \in \{0, \dots, h\} \wedge v \in V(G)\}$ and an edge (c, d) iff $\exists v \in V(G): \pi_i(v) = c \wedge \pi_{i+1}(v) = d$.
 101 We associate each tree node with the set of vertices of G having that color.¹ Here, we assume that
 102 the initial coloring is uniform. If this is not the case, we add an artificial root node and connect it to
 103 the initial colors. Likewise we insert the coloring $\pi_0 = \{V(G)\}$ as first element in the sequence of
 104 vertex colorings. An example color hierarchy is given in Figure 1.

105 Using this color hierarchy we can derive multiple colorings on the vertices: Choosing exactly one
 106 color on every path from the leaves to the root (or only the root), always leads to a valid coloring. The
 107 finest coloring is induced by the colors representing the leaves of the tree. Given a color hierarchy T ,
 108 we denote this coloring (which is equal to π_h) by π_T .

Weisfeiler-Leman Color Refinement. The 1-dimensional Weisfeiler-Leman (WL) algorithm or
 color refinement [16, 17] starts with a coloring c_0 , where all vertices have a color representing their
 label (or a uniform coloring in case of unlabeled vertices). In iteration i , the coloring c_i is obtained
 by assigning each vertex v in $V(G)$ a new color according to the colors of its neighbors, i.e.,

$$c_{i+1}(v) = z(c_i(v), \{\{c_i(u) \mid u \in N(v)\}\}),$$

109 where $z: \mathbb{N}_0 \times \mathbb{N}_0^{\mathbb{N}_0} \rightarrow \mathbb{N}_0$ is an injective function. Figure 1 depicts the first iterations of the algorithm
 110 for an example graph.

111 After enough iterations the number of different colors will no longer change and this resulting coloring
 112 is called the *coarsest stable coloring*. The coarsest stable coloring is unique and always reached after
 113 at most $|V(G)| - 1$ iterations. This trivial upper bound on the number of iterations is tight [18]. In
 114 practice, however, Weisfeiler-Leman refinement converges much faster (see Appendix C).

115 **Graph Kernels and the Weisfeiler-Leman Subtree Kernel.** A *kernel* on X is a function $k: X \times$
 116 $X \rightarrow \mathbb{R}$, so that there exist a Hilbert space \mathcal{H} and a mapping $\phi: X \rightarrow \mathcal{H}$ with $k(x, y) = \langle \phi(x), \phi(y) \rangle$
 117 for all x, y in X , where $\langle \cdot, \cdot \rangle$ is the inner product of \mathcal{H} . A *graph kernel* is a kernel on graphs, i.e., X
 118 is the set of all graphs.

119 The Weisfeiler-Leman subtree kernel [2] with height h is defined as

$$k_{ST}^h(G_1, G_2) = \sum_{i=0}^h \sum_{u \in V(G_1)} \sum_{v \in V(G_2)} \delta(c_i(u), c_i(v)), \quad (1)$$

120 where δ is the Dirac kernel (1, iff $c_i(u)$ and $c_i(v)$ are equal, and 0 otherwise). It counts the number of
 121 vertices with common colors in the two graphs up to the given bound on the number of Weisfeiler-
 122 Leman iterations.

¹Here, we consider a color hierarchy for a single graph, but given corresponding colorings on a set of graphs, a color hierarchy can be generated in the same way. In an inductive setting, where no fixed set of graphs is given, the color hierarchy contains all colors, that can be produced by the corresponding coloring algorithm.

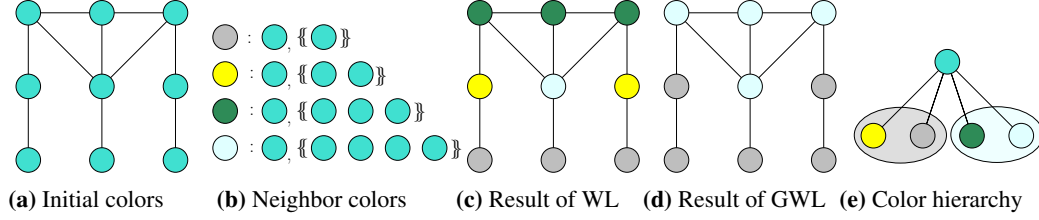


Figure 2: Initial coloring and results of the first iteration using WL and GWL refinement. We assume that the update function of GWL is a clustering algorithm, producing two clusters per old color. Vertices colored gray and yellow by WL are put into the same cluster, as well as green and light blue ones, as their neighbor color multisets only differ by one element each.

123 4 Gradual Weisfeiler-Leman Refinement

124 As a different approach to the refinement step of the Weisfeiler-Leman algorithm, we essentially
 125 replace the injective relabeling function with a non-injective one. We do this by allowing vertices
 126 with differing neighbor color multisets to be assigned the same color under some conditions. Through
 127 this, the number of colors per iteration can be limited, allowing to obtain a more gradual refinement
 128 of colors. To reach the same stable coloring as the original Weisfeiler-Leman algorithm, the function
 129 has to assure that vertices with differing colors in one iteration will get differing colors in future
 130 iterations and that in each iteration at least one color is split up, if possible.

131 We first define the property necessary to reach the stable coloring of the original Weisfeiler-Leman
 132 algorithm and discuss connections to the original as well as other vertex refinement algorithms. Then
 133 we provide a realization of such a function by means of clustering, which is used in our experimental
 134 evaluation. Figure 2 illustrates our idea. It depicts the initial coloring, the result of the first iteration
 135 of WL and a possible result of the first iteration of the gradual Weisfeiler-Leman refinement (GWL),
 136 when restricting the maximum number of new colors to 2 by clustering the neighbor color multisets.

Update Functions. Using the same approach as the Weisfeiler-Leman algorithm, the color of a vertex is updated iteratively according to the colors of its neighbors. Let \mathcal{T}_i denote a color hierarchy belonging to G and $n_i(v) = \{\{\pi_{\mathcal{T}_i}(x) \mid x \in N(v)\}\}$ the neighbor color multiset of v in iteration i . We use a similar update strategy, but generalize it using a special type of function:

$$\forall v \in V(G): c_{i+1}(v) = \pi_{\mathcal{T}_{i+1}}(v), \text{ with } \mathcal{T}_{i+1} = f(G, \mathcal{T}_i),$$

137 where f is a refining, neighborhood preserving function.

138 A refining, neighborhood preserving (*relep*) function f maps a pair (G, \mathcal{T}_i) to a tree \mathcal{T}_{i+1} , such that

139 **Condition 1.** $\mathcal{T}_i \subseteq \mathcal{T}_{i+1}$

140 **Condition 2.** $\mathcal{T}_i = \mathcal{T}_{i+1}$, iff $\forall v, w \in V(G) : v \approx_{\pi_{\mathcal{T}_i}} w \Rightarrow n_i(v) = n_i(w)$

141 **Condition 3.** $\mathcal{T}_i \subsetneq \mathcal{T}_{i+1} \Rightarrow \pi_{\mathcal{T}_{i+1}} \prec \pi_{\mathcal{T}_i}$

142 **Condition 4.** $\forall v, w \in V(G) : (v \approx_{\pi_{\mathcal{T}_i}} w \wedge n_i(v) = n_i(w)) \Rightarrow v \approx_{\pi_{\mathcal{T}_{i+1}}} w$

143 The conditions assure, that the coloring $\pi_{\mathcal{T}_{i+1}}$ is a strict refinement of $\pi_{\mathcal{T}_i}$, if there exists a strict
 144 refinement: Condition 1 assures that the new coloring is a refinement of the old one. Condition 2
 145 assures that the tree (and in turn the coloring) only stays the same, iff the stable coloring is reached,
 146 while Condition 3 assures that, if the trees are not equal, $\pi_{\mathcal{T}_{i+1}}$ is a strict refinement of $\pi_{\mathcal{T}_i}$. Without
 147 this condition it would be possible to obtain a tree, that fulfills Condition 1 but does not strictly refine
 148 the coloring (for example by adding one child to each leaf). Condition 4 assures that vertices, that are
 149 indistinguishable regarding their color and their neighbor color multiset, get the same color (as in the
 150 original Weisfeiler-Leman algorithm).

151 We call this new approach *gradual Weisfeiler-Leman refinement* (GWL refinement). Since f is a
 152 *relep* function, it is assured that at least one color is split into at least two new colors, if the stable
 153 coloring is not yet reached. This property and its implications are explored in the following section.

154 Usually, the refinement is computed simultaneously for multiple graphs. This can be realized by using
 155 the disjoint union of all graphs as input. Note that this will have an influence on the function f , since
 156 refinements might differ based on the vertices involved. This is a typical case of transductive learning,

157 because the algorithm has to run on all graphs and if a new graph is encountered, the algorithm has to
 158 run again on the enlarged graph. To counteract this, one could only run the algorithm on the training
 159 set and keep track of the coloring choices. Then for new graphs, the vertices are colored with the
 160 most similar colors.

161 4.1 Equivalence of the Stable Colorings

162 The gradual color refinement will never assign two vertices the same color, if their colors differed in
 163 the previous iteration, since we require the coloring to be a refinement of the previous one. We can
 164 show that the stable coloring obtained by GWL refinement using any renep function is equal to the
 165 unique coarsest stable coloring, which is obtained by the original Weisfeiler-Leman algorithm.

166 **Theorem 5** ([19], Proposition 3). *For every coloring π of $V(G)$, there is a unique coarsest stable*
 167 *coloring p that refines π .*

168 This means GWL with any renep function, should it reach a coarsest stable coloring, will reach this
 169 unique coarsest stable coloring. It remains to show that GWL will reach a coarsest stable coloring.

170 **Theorem 6.** *For all G the GWL refinement using any renep function will find the unique coarsest*
 171 *stable coloring of $V(G)$.*

172 *Proof.* Let $\pi_{\mathcal{T}} = \{p_1, \dots, p_n\}$ be the stable coloring obtained from GWL on the initial coloring
 173 π_0 . Assume there exists another stable coloring $\pi' = \{p'_1, \dots, p'_m\}$ with $\pi_{\mathcal{T}} \prec \pi' \preceq \pi_0$, so $m < n$.
 174 Then $\exists v, w \in V(G) : (v \approx_{\pi'} w \wedge v \not\approx_{\pi_{\mathcal{T}}} w)$ and since Condition 4 applies $n(v) \neq n(w)$, which
 175 contradicts the assumption that π' is stable. \square

176 The original Weisfeiler-Leman refinement can be realized by using the renep function with \Leftrightarrow instead
 177 of \Rightarrow in Condition 4. This ensures that vertices get assigned the same color, iff they previously had
 178 the same color and their neighborhood color multisets do not differ. Since this procedure splits all
 179 colors, that can be split up, it is the fastest converging possible renep function (because only direct
 180 neighborhood is considered). A trivial upper bound for the maximum number of Weisfeiler-Leman
 181 iterations needed is $|V(G)| - 1$ and there are infinitely many graphs on which this number of iterations
 182 is required for convergence [18]. We obtain the same upper bound for GWL.

183 **Theorem 7.** *The maximum number of iterations needed to reach the stable coloring using GWL*
 184 *refinement is $|V(G)| - 1$.*

185 *Proof.* The function we consider is a renep function. It follows that, prior to reaching the stable
 186 coloring, at least one color is split into at least two new colors in every iteration. Since vertices that
 187 had different colors at any step will also have different colors in the following iterations, the number
 188 of colors increases in every step. Hence, after at most $|V(G)| - 1$ steps, each vertex has a unique
 189 color, which is a stable coloring. \square

190 **Sequential Weisfeiler-Leman.** For optimizing the running time of the Weisfeiler-Leman algorithm,
 191 sequential refinement strategies have been proposed [19–21], which lead to the same stable coloring
 192 as the original WL. Our presentation follows Berkholz et al. [19], who provide implementation
 193 details and a thorough complexity analysis. Sequential WL manages a stack containing the colors
 194 that still have to be processed. All initial colors are added to this stack. In each step, the next
 195 color c from the stack is used to refine the current coloring π (and generate a new coloring π')
 196 using the following update strategy: $\forall v, w \in V(G) : v \approx_{\pi'} w \Leftrightarrow \left| \left\{ \{x \mid x \in N(v) \wedge \pi(x) = c\} \right\} \right| =$
 197 $\left| \left\{ \{x \mid x \in N(w) \wedge \pi(x) = c\} \right\} \right| \wedge v \approx_{\pi} w$. Note that $\pi' \prec \pi$ is not guaranteed. For colors that are
 198 split, all new colors are added to the stack with exception of the largest color class. This is shown to
 199 be sufficient for generating the coarsest stable coloring [19].

200 Sequential Weisfeiler-Leman can be realized by our GWL with the restriction, that in sequential
 201 WL, some refinement operations might not produce strict refinements. We need to skip these in our
 202 approach (since renep functions have to produce strict refinements as long as the coloring is not stable).
 203 The renep function has to fulfill $\forall v, w \in V(G) : v \approx_{\pi_{\mathcal{T}_{i+1}}} w \Leftrightarrow \left| \left\{ \{x \mid x \in N(v) \wedge \pi_{\mathcal{T}_i}(x) = c\} \right\} \right| =$
 204 $\left| \left\{ \{x \mid x \in N(w) \wedge \pi_{\mathcal{T}_i}(x) = c\} \right\} \right| \wedge v \approx_{\pi_{\mathcal{T}_i}} w$, where c is the next color in the stack that produces a
 205 strict refinement.

206 4.2 Running Time

207 The running time of the gradual Weisfeiler-Leman refinement depends on the cost of the update
208 function used.

209 **Theorem 8.** *The running time for the gradual Weisfeiler-Leman refinement is $O(h \cdot t_u(|V(G)|))$,*
210 *where h is the number of iterations and $t_u(n)$ is the time needed to compute the renep function for n*
211 *elements.*

212 The update function used in the original Weisfeiler-Leman refinement can be computed in time
213 $O(|V(G)| + |E(G)|)$ in the worst-case by sorting the neighbor color multisets using bucket sort [2].

214 4.3 Discussion of Suitable Update Functions

215 The update function of the original Weisfeiler-Leman refinement provides a fast way to reach the
216 stable coloring, but in machine learning tasks a more fine grained vertex similarity is needed. A
217 suitable update function restricts the number of new colors to a manageable amount, while still
218 fulfilling the requirements of a renep function. Clustering the neighborhood multisets of the vertices,
219 and letting the clusters imply the new colors, is an intuitive way to restrict the number of colors per
220 iteration and assign similar neighborhoods the same new color. We discuss how to realize a renep
221 function using clustering.

222 Whether two vertices, that currently have the same color, will be assigned the same color in the next
223 step, depends on two factors: If they have the same neighbor color multiset, they have to remain in
224 one color group. If their neighbor color multisets differ, however, the renep function can decide to
225 either separate them or not (provided any new colors are generated to fulfill Condition 3). We propose
226 clustering the neighbor color multisets separately for each old color and let the clusters imply new
227 colors. If a clustering function guarantees to produce at least two clusters for inputs with at least two
228 distinct objects, we obtain a renep function.

229 Although various clustering algorithms are available, we identified k -means as a convenient choice
230 because of its efficiency and controllability of the number of clusters. In order to apply k -means to
231 multisets of colors, we represent them as (sparse) vectors, where each entry counts the number of
232 neighbors with a specific color. The above method using k -means clustering with $k > 1$ satisfies the
233 requirements of a renep function. Of course, if the number of elements to cluster is less than or equal
234 to k , the clustering can be omitted and each element can be assigned its own cluster. The number of
235 clusters in iteration i is bounded by $|L| \cdot k^i$, since each color can split into at most k new colors in
236 each iteration and the initial coloring has at most $|L|$ colors.

237 5 Applications

238 The gradual Weisfeiler-Leman refinement provides a more fine-grained approach to capture vertex
239 similarity, where two vertices are considered more similar, the longer it takes until they get assigned
240 different colors. This makes the approach applicable not only to vertex classification, but also in graph
241 kernels and as a vertex similarity measure for graph matching. We further describe these possible
242 applications in the following and evaluate them against the state-of-the-art methods in Section 6.

243 **Graph Kernels.** The idea of the GWL subtree kernel is essentially the same as the Weisfeiler-
244 Leman subtree kernel [2], but instead of using the original Weisfeiler-Leman algorithm, the GWL
245 algorithm is used to generate the features. We use the definition given in Equation (1) replacing
246 the Weisfeiler-Leman colorings with the coloring from the GWL algorithm. The Weisfeiler-Leman
247 optimal assignment kernel [5] is obtained from an optimal assignment between the vertices of two
248 graphs regarding a vertex similarity obtained from a color hierarchy. We replace the Weisfeiler-Leman
249 color hierarchy used originally by the one from our gradual refinement. We evaluate the performance
250 of our newly proposed kernels in Section 6.

251 **Tree Metrics for Approximating the Graph Edit Distance.** The graph edit distance, a general
252 distance measurement for graphs, is usually approximated due to its complexity. Many approxima-
253 tions find an optimal assignment between the vertices of the two graphs and derive a (sub-optimal)
254 edit path from this assignment. Naturally, the cost of such an edit path is an upper bound for the
255 graph edit distance. An optimal assignment can be computed in linear time, if the underlying cost

function is a tree metric [22], which means an approximation of the graph edit distance can be found in linear time, given a tree metric on the vertices. The color hierarchy, see Figure 1, produced by the original Weisfeiler-Leman refinement, as well as our gradual variant, can be interpreted as such a tree metric. Lin [22] approximated the graph edit distance as described above. We can again replace the original color hierarchy by the one produced by our gradual Weisfeiler-Leman refinement. Appendix I includes a more detailed explanation on how to approximate the graph edit distance using assignments. We evaluate this approximation of the graph edit distance regarding its accuracy in k -nn-classification against the state-of-the-art and the original approach.

6 Experimental Evaluation

We evaluate the proposed approach regarding its applicability in graph kernels, as the gradual Weisfeiler-Leman subtree kernel (GWL) and the gradual Weisfeiler-Leman optimal assignment kernel (GWLOA), as well as its usefulness as a tree metric for approximating the graph edit distance. Specifically, we address the following research questions:

- Q1** Can our kernels compete with state-of-the-art methods regarding classification accuracy on real-world and synthetic datasets?
- Q2** Which refinement speed is appropriate and are there dataset-specific differences?
- Q3** How do our kernels compare to the state-of-the-art methods in terms of running time?
- Q4** Is the vertex similarity obtained from GWL refinement suitable for learning with approximated graph edit distance?

We compare to the Weisfeiler-Leman subtree kernel (WLST) [2], the Weisfeiler-Leman optimal assignment kernel (WLOA) [5], as well as RWL* [14], the approximation of the relaxed Weisfeiler-Leman subtree kernel, and the deep Weisfeiler-Leman kernel (DWL) [12]. We do not compare to [4], since the kernel showed results similar to the WLOA kernel. We compare the graph edit distance approximation using our tree metric GWLT to the original approach Lin [22] and state-of-the-art method BGM [23].

6.1 Setup

As discussed in Section 4.3 we used k -means clustering in our new approach. If for any color less than k different vectors were present in the clustering step, each distinct vector got its own cluster. We implemented our GWL, GWLOA and also the original WLST and WLOA in Java. We used the RWL* and DWL Python implementations provided by the authors. Note that in contrast to the other approaches, this RWL* uses multi-threading.

For evaluation, we used the C -SVM implementation LIBSVM [24] and report average classification accuracies obtained by 10-fold nested cross-validation repeated 10 times with random fold assignments. The parameters of the SVM and the kernels were optimized by the inner cross-validation on the current training fold using grid search. We chose $C \in \{10^{-3}, 10^{-2}, \dots, 10^3\}$ for the SVM, $h \in \{0, \dots, 10\}$ for WLST, WLOA, GWL and GWLOA and k -means with $k \in \{2, 4, 8, 16\}$. For RWL* we used $h \in \{1, \dots, 4\}$ and default values for the other parameters. In DWL the window size w and dimension d were set to 25, since this generally worked best out of the combinations from $d, w \in \{5, 25, 50\}$ and no defaults were provided. We used the default settings for the other parameters and again $h \in \{1, \dots, 10\}$. The running time experiments were conducted on an Intel Xeon Gold 6130 machine at 2.1 GHz with 96 GB RAM. For evaluating the approximation of the graph edit distance, we compare the 1- nn classification accuracy and used the Java implementation of Lin provided by the authors and implemented our approach GWLT, as well as BGM, also in Java for a fair comparison.

Extension to Edge Labels. The original Weisfeiler-Leman algorithm can be extended to respect edge labels by updating the colors according to $c_{i+1}(v) = z(c_i(v), \{\{(l(u, v), c_i(u)) \mid u \in N(v)\}\})$. All kernels used in the comparison use a similar strategy to incorporate edge labels if present.

Datasets. We used several real-world datasets from the TUDataset [25] and the *EGO-Nets* datasets [14] for our experiments. See Appendix B and F for an overview of the datasets, as well as additional synthetic datasets and corresponding results. We selected these datasets as they

Table 1: Average classification accuracy and standard deviation (highest accuracies marked in **bold**).

Kernel	PTC_FM	KKI	EGO-1	EGO-2	EGO-3	EGO-4
WLST	64.16 \pm 1.30	49.97 \pm 2.88	51.30 \pm 2.42	57.15 \pm 1.61	56.15 \pm 1.67	53.40 \pm 1.77
DWL	64.18 \pm 1.46	50.93 \pm 2.87	55.80 \pm 1.35	56.50 \pm 1.64	55.90 \pm 1.64	53.25 \pm 2.81
RWL*	62.43 \pm 1.46	46.54 \pm 4.03	65.60 \pm 2.74	70.20 \pm 1.36	67.60 \pm 1.07	74.25 \pm 2.12
WLOA	62.34 \pm 1.39	48.72 \pm 4.05	55.95 \pm 1.11	60.30 \pm 2.00	54.25 \pm 1.35	52.30 \pm 2.29
GWL	62.61 \pm 1.94	57.79 \pm 3.95	67.95 \pm 2.05	73.65 \pm 1.86	65.45 \pm 1.88	77.45 \pm 1.97
GWLOA	64.58 \pm 1.77	47.47 \pm 2.41	69.80 \pm 1.65	72.40 \pm 2.52	67.45 \pm 1.69	75.35 \pm 1.67
	COLLAB	DD	IMDB-B	MSRC_9	NCI1	REDDIT-B
WLST	78.98 \pm 0.22	79.00 \pm 0.52	72.01 \pm 0.80	90.13 \pm 0.75	85.96 \pm 0.18	80.81 \pm 0.52
DWL	78.93 \pm 0.18	78.92 \pm 0.40	72.36 \pm 0.56	90.50 \pm 0.76	85.68 \pm 0.18	80.83 \pm 0.40
RWL*	77.94 \pm 0.38	77.52 \pm 0.65	72.96 \pm 0.86	88.86 \pm 0.89	79.45 \pm 0.32	77.69 \pm 0.31
WLOA	80.81 \pm 0.22	79.44 \pm 0.31	72.60 \pm 0.89	90.68 \pm 0.92	86.29 \pm 0.13	89.40 \pm 0.14
GWL	80.62 \pm 0.33	79.00 \pm 0.81	73.66 \pm 1.25	88.32 \pm 1.20	85.33 \pm 0.35	86.46 \pm 0.35
GWLOA	81.30 \pm 0.29	78.49 \pm 0.57	72.88 \pm 0.79	91.27 \pm 1.06	85.36 \pm 0.36	89.98 \pm 0.34

306 cover a wide range of applications, consisting of both molecule datasets and graphs derived from
 307 social networks. See Appendix C for the number of Weisfeiler-Leman iterations needed to reach the
 308 stable coloring for each dataset.

309 6.2 Results

310 In the following, we present the classification accuracy, as well as running time, of the different
 311 kernel methods. We investigate the parameter selection for our algorithm and discuss the application
 312 of our approach for approximating the graph edit distance.

313 **Q1: Classification Accuracy.** Table 1 shows the classification accuracy of the different kernels.
 314 While on some datasets our new approaches do not outcompete all state-of-the-art methods, they
 315 are more accurate in most cases, in some cases even with a large margin to the second-placed (for
 316 example on *KKI*, *EGO-1* or *EGO-4*). While *RWL** is better than our approaches on some datasets,
 317 the running time of this method is much higher, cf. Q3. *WLOA* also produces very good results on
 318 many datasets, but cannot compete on the *EGO-Nets* and synthetic datasets (see Appendix F). For
 319 molecular graphs (*PTC_FM*, *NCI*) we see no significant improvements, which can be explained by
 320 their small degree and sensitivity of molecular properties to small changes. Overall, our method
 321 provides the highest accuracy on 9 of 12 datasets and is close to the best accuracy for the others.

322 **Q2: Parameter Selection.** For *GWL* and *GWLOA* two parameters have to be chosen: The number
 323 of iterations h and the number k of clusters in k -means. We investigate which choices lead to the best
 324 classification accuracy. Figure 3 shows the number of times, a specific parameter combination was
 325 selected as it provided the best accuracy for the test set. Here, we only show the parameter selection
 326 for some of the datasets. The results for the other datasets, as well as the parameter selection for
 327 *WLST* and *WLOA*, can be found in Appendix D. We can see that for *GWL* and most datasets the
 328 best k is in $\{2, 4, 8\}$ and on those datasets classification accuracy of *GWL* exceeds that of *WLST*. On
 329 datasets on which *GWL* performed worse than *WLST*, the best choice for parameters is not clear
 330 and it seems like a larger k might be beneficial for improving the classification accuracy. Similar
 331 tendencies can be observed for *GWLOA*.

332 **Q3: Running Time.** Figure 4 shows the time needed for computing the feature vectors using the
 333 different kernels (for results on the other datasets and the influence of the parameter k on running
 334 time see Appendix E and G). *RWL** and *DWL* are much slower than the other kernels, while only
 335 *RWL** leads to minor improvements in classification accuracy on few datasets. While our approach is
 336 only slightly slower than *WLST/WLOA*, it yields great improvements on the classification accuracy
 337 on most datasets, cf. Q1.

338 **Q4: Learning with Approximated Graph Edit Distance.** Table 2 compares the 1-nn classifi-
 339 cation accuracy of our approach when approximating the graph edit distance to the original [22]

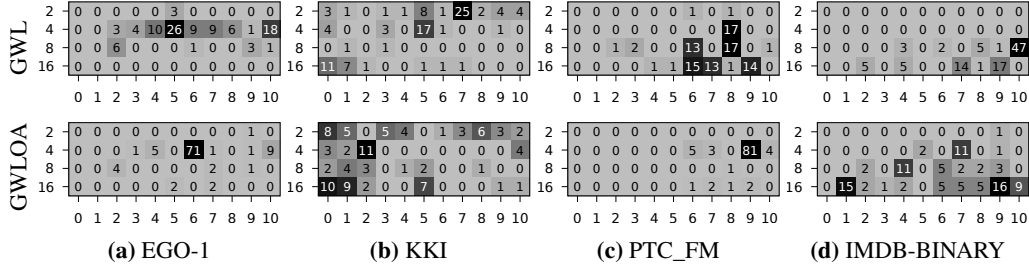


Figure 3: Number of times a parameter combination of GWL and GWLOA was selected from $k \in \{2, 4, 8, 16\}$ and $h \in \{0, \dots, 10\}$ based on the accuracy achieved on the test set.

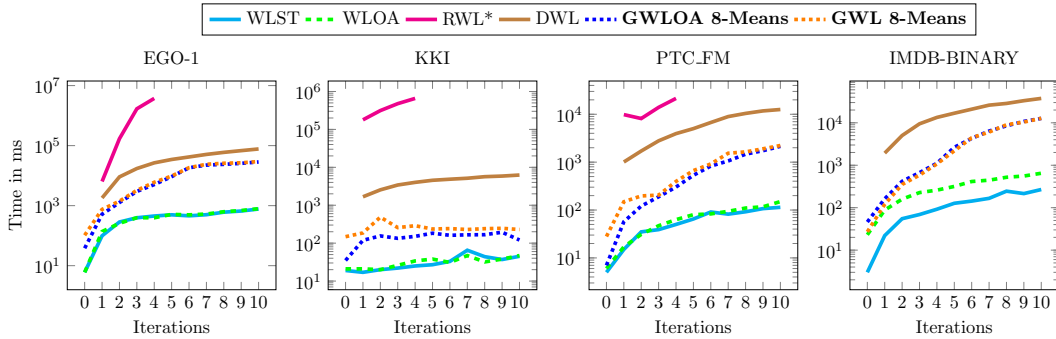


Figure 4: Running time in milliseconds for computing the feature vectors for all graphs of a dataset using the different methods. Note that RWL* uses multi-threading, while the other methods do not. Missing values for RWL* and DWL in the larger datasets are due to timeout.

340 and another state-of-the-art method based on bipartite graph matching [23]. Our approach clearly
 341 outcompetes both methods on all datasets. In Appendix J we investigate the approximation quality
 342 for similar graphs (which are important to k -nn classification) further.

343 7 Conclusions

344 We proposed a general framework for iterative vertex refinement generalizing the popular Weisfeiler-
 345 Leman algorithm and discussed connections to other vertex refinement strategies. Based on this, we
 346 proposed two new graph kernels and showed that they outperform the original Weisfeiler-Leman
 347 subtree kernel and similar state-of-the-art approaches in terms of classification accuracy in almost all
 348 cases, while keeping the running time much lower than comparable methods. We also investigated
 349 the application of our method to approximating the graph edit distance, where we again outperformed
 350 the state-of-the-art methods.

351 In further research it might be interesting to systematically compare our approach to graph neural
 352 networks, since their message passing scheme is similar to the update strategy of the WL algorithm.
 353 Moreover, other renep functions can be explored, for example, by using other clustering strategies, or
 354 by developing new concepts for inexact neighborhood comparison.

Table 2: Average classification accuracy and standard deviation (highest accuracies marked in **bold**).

Method	PTC_FM	MSRC_9	KKI	EGO-1	EGO-2	EGO-3	EGO-4
BGM	60.14 \pm 1.50	72.13 \pm 1.28	43.89 \pm 1.27	44.75 \pm 1.05	42.05 \pm 1.25	out of time	out of time
Lin	62.38 \pm 1.08	81.36 \pm 0.64	55.18 \pm 2.44	40.40 \pm 1.17	31.65 \pm 1.07	26.60 \pm 0.94	36.55 \pm 1.72
GWLT	63.19 \pm 0.11	85.97 \pm 0.59	55.18 \pm 2.44	56.20 \pm 1.42	47.90 \pm 1.28	36.40 \pm 1.04	47.90 \pm 1.32

References

- 355
- 356 [1] Brendan D. McKay and Adolfo Piperno. Practical graph isomorphism, II. *J. Symb. Comput.*,
357 60:94–112, 2014. doi: 10.1016/j.jsc.2013.09.003. 1
- 358 [2] Nino Shervashidze, Pascal Schweitzer, Erik Jan van Leeuwen, Kurt Mehlhorn, and Karsten M.
359 Borgwardt. Weisfeiler-Lehman graph kernels. *J. Mach. Learn. Res.*, 12:2539–2561, 2011. URL
360 <http://dl.acm.org/citation.cfm?id=2078187>. 1, 2, 3, 6, 7
- 361 [3] Matteo Togninalli, M. Elisabetta Ghisu, Felipe Llinares-López, Bastian Rieck, and Karsten M.
362 Borgwardt. Wasserstein Weisfeiler-Lehman graph kernels. In Hanna M. Wallach, Hugo
363 Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett,
364 editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural
365 Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC,
366 Canada*, pages 6436–6446, 2019. 1
- 367 [4] Dai Hai Nguyen, Canh Hao Nguyen, and Hiroshi Mamitsuka. Learning subtree pattern im-
368 portance for Weisfeiler-Lehman based graph kernels. *Mach. Learn.*, 110(7):1585–1607, 2021.
369 7
- 370 [5] Nils M. Kriege, Pierre-Louis Giscard, and Richard C. Wilson. On valid optimal assignment
371 kernels and applications to graph classification. In *Advances in Neural Information Processing
372 Systems*, pages 1615–1623, 2016. 1, 2, 6, 7
- 373 [6] Christopher Morris, Matthias Fey, and Nils M. Kriege. The power of the Weisfeiler-Leman
374 algorithm for machine learning with graphs. In Zhi-Hua Zhou, editor, *Proceedings of the
375 Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event /
376 Montreal, Canada, 19-27 August 2021*, pages 4543–4550. ijcai.org, 2021. doi: 10.24963/ijcai.
377 2021/618. 1
- 378 [7] Christopher Morris, Yaron Lipman, Haggai Maron, Bastian Rieck, Nils M. Kriege, Martin
379 Grohe, Matthias Fey, and Karsten M. Borgwardt. Weisfeiler and Leman go machine learning:
380 The story so far. *CoRR*, abs/2112.09992, 2021. 1
- 381 [8] Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen,
382 Gaurav Rattan, and Martin Grohe. Weisfeiler and Leman go neural: Higher-order graph neural
383 networks. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, Honolulu,
384 Hawaii, USA, January 27 - February 1, 2019*, pages 4602–4609. AAAI Press, 2019. doi:
385 10.1609/aaai.v33i01.33014602. 1
- 386 [9] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural net-
387 works? In *7th International Conference on Learning Representations, ICLR*. OpenReview.net,
388 2019. 1
- 389 [10] Nils M. Kriege, Fredrik D. Johansson, and Christopher Morris. A survey on graph kernels.
390 *Applied Network Science*, 5(1):6, 2020. doi: 10.1007/s41109-019-0195-3. 2
- 391 [11] Karsten M. Borgwardt, M. Elisabetta Ghisu, Felipe Llinares-López, Leslie O’Bray, and Bastian
392 Rieck. Graph kernels: State-of-the-art and future challenges. *Found. Trends Mach. Learn.*, 13
393 (5-6), 2020. doi: 10.1561/22000000076. 2
- 394 [12] Pinar Yanardag and S. V. N. Vishwanathan. Deep graph kernels. In Longbing Cao, Chengqi
395 Zhang, Thorsten Joachims, Geoffrey I. Webb, Dragos D. Margineantu, and Graham Williams,
396 editors, *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Dis-
397 covery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, pages 1365–1374. ACM,
398 2015. doi: 10.1145/2783258.2783417. 2, 7
- 399 [13] Annamalai Narayanan, Mahinthan Chandramohan, Lihui Chen, Yang Liu, and Santhoshkumar
400 Saminathan. subgraph2vec: Learning distributed representations of rooted sub-graphs from
401 large graphs. *CoRR*, abs/1606.08928, 2016. URL <http://arxiv.org/abs/1606.08928>. 2
- 402 [14] Till Hendrik Schulz, Tamás Horváth, Pascal Welke, and Stefan Wrobel. A generalized
403 Weisfeiler-Lehman graph kernel. *Mach Learn*, 2022. 2, 7, 12
- 404 [15] Zhongxi Fang, Jianming Huang, Xun Su, and Hiroyuki Kasai. Wasserstein graph distance
405 based on l_1 -approximated tree edit distance between Weisfeiler-Lehman subtrees. 2022. doi:
406 10.48550/ARXIV.2207.04216. URL <https://arxiv.org/abs/2207.04216>. 2

- 407 [16] Martin Grohe, Kristian Kersting, Martin Mladenov, and Pascal Schweitzer. Color Refinement
408 and Its Applications. In *An Introduction to Lifted Probabilistic Inference*. The MIT Press, 08
409 2021. ISBN 9780262365598. 3
- 410 [17] Boris Weisfeiler and Adrian Leman. A reduction of a graph to canonical form and an algebra
411 arising during this reduction. *Nauchno-Tekhnicheskaya Informatsia*, 2(9):12–16, 1968. 3
- 412 [18] Sandra Kiefer and Brendan D. McKay. The iteration number of colour refinement. In Ar-
413 tur Czumaj, Anuj Dawar, and Emanuela Merelli, editors, *47th International Colloquium on*
414 *Automata, Languages, and Programming, ICALP 2020, July 8-11, 2020, Saarbrücken, Ger-*
415 *many (Virtual Conference)*, volume 168 of *LIPIcs*, pages 73:1–73:19. Schloss Dagstuhl -
416 Leibniz-Zentrum für Informatik, 2020. doi: 10.4230/LIPIcs.ICALP.2020.73. URL <https://doi.org/10.4230/LIPIcs.ICALP.2020.73>. 3, 5
- 418 [19] Christoph Berkholz, Paul S. Bonsma, and Martin Grohe. Tight lower and upper bounds for the
419 complexity of canonical colour refinement. *Theory Comput. Syst.*, 60(4):581–614, 2017. doi:
420 10.1007/s00224-016-9686-0. URL <https://doi.org/10.1007/s00224-016-9686-0>. 5
- 421 [20] Robert Paige and Robert Endre Tarjan. Three partition refinement algorithms. *SIAM J. Comput.*,
422 16(6):973–989, 1987. doi: 10.1137/0216062. URL <https://doi.org/10.1137/0216062>.
- 423 [21] Tommi A. Junttila and Petteri Kaski. Engineering an efficient canonical labeling tool for
424 large and sparse graphs. In *Proceedings of the Nine Workshop on Algorithm Engineering and*
425 *Experiments, ALENEX 2007, New Orleans, Louisiana, USA, January 6, 2007*. SIAM, 2007. doi:
426 10.1137/1.9781611972870.13. URL <https://doi.org/10.1137/1.9781611972870.13>. 5
- 427 [22] Nils M. Kriege, Pierre-Louis Giscard, Franka Bause, and Richard C. Wilson. Computing
428 optimal assignments in linear time for approximate graph matching. In *ICDM*, pages 349–358,
429 2019. 7, 8, 16, 17
- 430 [23] Kaspar Riesen and Horst Bunke. Approximate graph edit distance computation by means of
431 bipartite graph matching. *Image Vision Comput.*, 27(7):950–959, 2009. 7, 9, 16
- 432 [24] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM*
433 *Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at
434 <http://www.csie.ntu.edu.tw/~cjlin/libsvm>. 7
- 435 [25] Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion
436 Neumann. TUDataset: A collection of benchmark datasets for learning with graphs. In
437 *ICML 2020 Workshop on Graph Representation Learning and Beyond, GRL+*, 2020. URL
438 <http://www.graphlearning.io>. 7, 12
- 439 [26] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. LIB-
440 LINEAR: A library for large linear classification. *Journal of Machine Learning Research*, 9:
441 1871–1874, 2008. 16
- 442 [27] Zhiping Zeng, Anthony K. H. Tung, Jianyong Wang, Jianhua Feng, and Lizhu Zhou. Comparing
443 stars: On approximating graph edit distance. *Proc. VLDB Endow.*, 2(1):25–36, 2009. 16

Table 3: Datasets with discrete vertex and edge labels and their statistics [25]. The *EGO-Nets* datasets [14] are unlabeled.

Name	Graphs	Classes	avg V	avg E	L _V	L _E
<i>KKI</i>	83	2	26.96	48.42	190	–
<i>PTC_FM</i>	349	2	14.11	14.48	18	4
<i>COLLAB</i>	5000	3	74.49	2457.78	–	–
<i>DD</i>	1178	2	284.32	715.66	82	–
<i>IMDB-BINARY</i>	1000	2	19.77	96.53	–	–
<i>MSRC_9</i>	221	2	40.58	97.94	10	–
<i>NCII</i>	4110	2	29.87	32.30	37	–
<i>REDDIT-BINARY</i>	2000	2	429.63	497.75	–	–
<i>EGO-1</i>	200	4	138.97	593.53	–	–
<i>EGO-2</i>	200	4	178.55	1444.86	–	–
<i>EGO-3</i>	200	4	220.01	2613.49	–	–
<i>EGO-4</i>	200	4	259.78	4135.80	–	–

444 A Pseudocode

445 Algorithm 1 shows the procedure of gradual Weisfeiler-Leman refinement. We start with an initial
446 coloring (either uniform or based on the vertex labels) and then iteratively refine these colors using a
renew function, h times.

Algorithm 1 Gradual Weisfeiler-Leman Refinement.

```

1: procedure GWL( $G, h$ )                                     ▷  $h$  is number of iterations
2:   for all  $v \in V(G)$  do                                     ▷ Initial coloring
3:      $c_0(v) \leftarrow \mu(v)$ 
4:   initialize  $T_0$  as described in Section 3
5:   for  $i \leftarrow 1, i \leq h, i \leftarrow i + 1$  do
6:      $\mathcal{T}_i \leftarrow f(G, \mathcal{T}_{i-1})$                          ▷ Compute new colors
7:     for all  $v \in V(G)$  do                                   ▷ Assign new colors
8:        $c_i(v) \leftarrow \pi_{\mathcal{T}_i}(v)$ 

```

447

448 B Datasets

449 We used several real-world datasets from the TUDataset [25], the *EGO-Nets* datasets [14], as well
450 as synthetic datasets for our experiments. See Table 3 for an overview of the real-world datasets.
451 We selected these datasets as they cover a wide range of applications, consisting of both molecule
452 datasets and graphs derived from social networks.

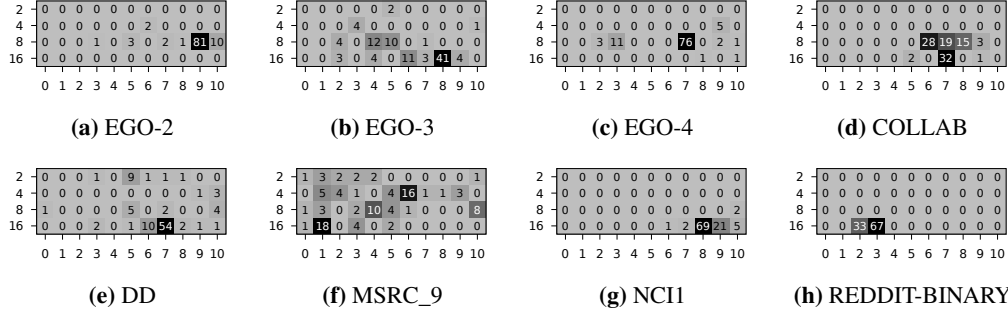
453 The synthetic datasets were generated using the block graph generation method [14]. We generated 9
454 synthetic datasets with two classes and 200 graphs in each class. For each dataset we first generated
455 two seed graphs (one per class) with 16 vertices, that both are constructed from a tree by appending
456 a single edge, so that their sets of vertex degrees are equal. For the dataset graphs each vertex of
457 the seed graph was replaced by 8 vertices. Vertices generated from the same seed vertex, as well as
458 from adjacent vertices are connected with probability p . m noise edges are then added randomly.
459 We investigated the two cases $p = 1.0$ and $m \in \{0, 10, 20, 50, 100\}$, and $p \in \{1.0, 0.8, 0.6, 0.4, 0.2\}$
460 and $m = 0$. We denote the datasets by S_{p_m} .

461 C Number of Weisfeiler-Leman Iterations

462 We investigate the number of WL iterations needed to reach the stable coloring in the various datasets.
463 On the datasets with – entries (see Table 4) our algorithm, that checked whether the stable coloring
464 is reached, did not finish in a reasonable time due to the size of the datasets/graphs. It can be seen

Table 4: Number of iterations needed to reach the stable coloring.

Dataset	<i>KKI</i>	<i>PTC_FM</i>	<i>COLLAB</i>	<i>DD</i>	<i>IMDB-B</i>	<i>MSRC_9</i>
WL	3	13	—	—	3	3
Dataset	<i>NCII</i>	<i>REDDIT-B</i>	<i>EGO-1</i>	<i>EGO-2</i>	<i>EGO-3</i>	<i>EGO-4</i>
WL	39	—	5	4	4	5

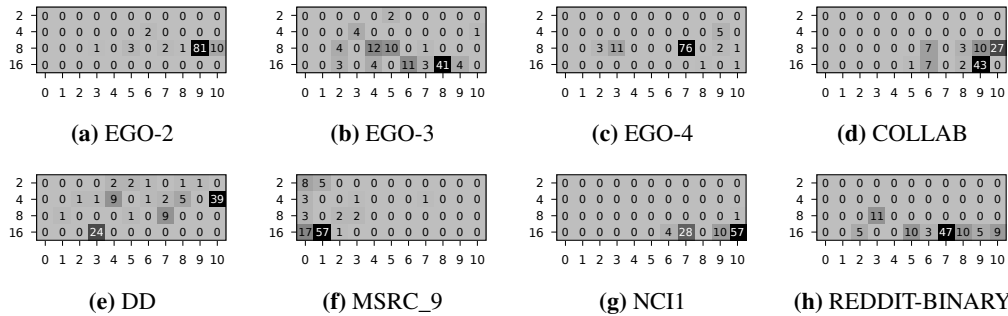

Figure 5: With $k \in \{2, 4, 8, 16\}$ and $h \in \{0, \dots, 10\}$, we show the number of times a specific parameter combination for GWL was selected as it provided the best accuracy for the test set.

465 that on most datasets the number of iterations needed to reach the stable coloring is very low. This
 466 means that after a few iterations we do not gain any new information when using for example the
 467 traditional WLST kernel.

468 D Parameter Selection - Further Results

469 Figures 5 and 6 show the parameter selection for the remaining datasets for GWL and GWLOA. In
 470 most datasets the choice is restricted to two or three values. For some of the datasets, the best choice
 471 seems to include $k = 16$. This might indicate that a larger k could be beneficial for increasing the
 472 accuracy.

473 Figure 7 shows the parameter selection for WLST and WLOA. There is only one parameter, the
 474 number of WL iterations, for both kernels. We can see that indeed on most datasets, only few
 475 iterations are needed to gain the best possible accuracy. On datasets such as *EGO-2*, *EGO-4* or *NCII*,
 476 however, this is not the case. For *NCII* we can assume that we still gain information through more
 477 iterations, since we have not yet reached the stable coloring. For the *EGO*-datasets, this is surprising,
 478 since the stable coloring is reached after 4 (5) iterations. On the other hand, the classification accuracy
 479 reached is still not good.


Figure 6: With $k \in \{2, 4, 8, 16\}$ and $h \in \{0, \dots, 10\}$, we show the number of times a specific parameter combination for GWLOA was selected as it provided the best accuracy for the test set.

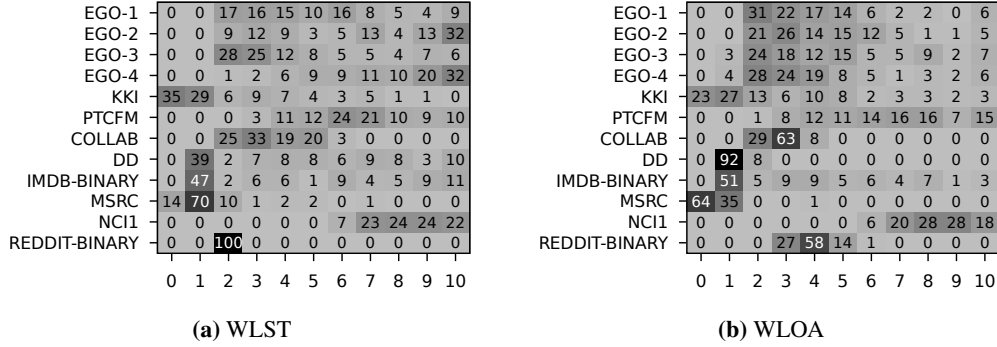


Figure 7: With $h \in \{0, \dots, 10\}$, we show the number of times a specific parameter combination for WLST and WLOA was selected as it provided the best accuracy for the test set.

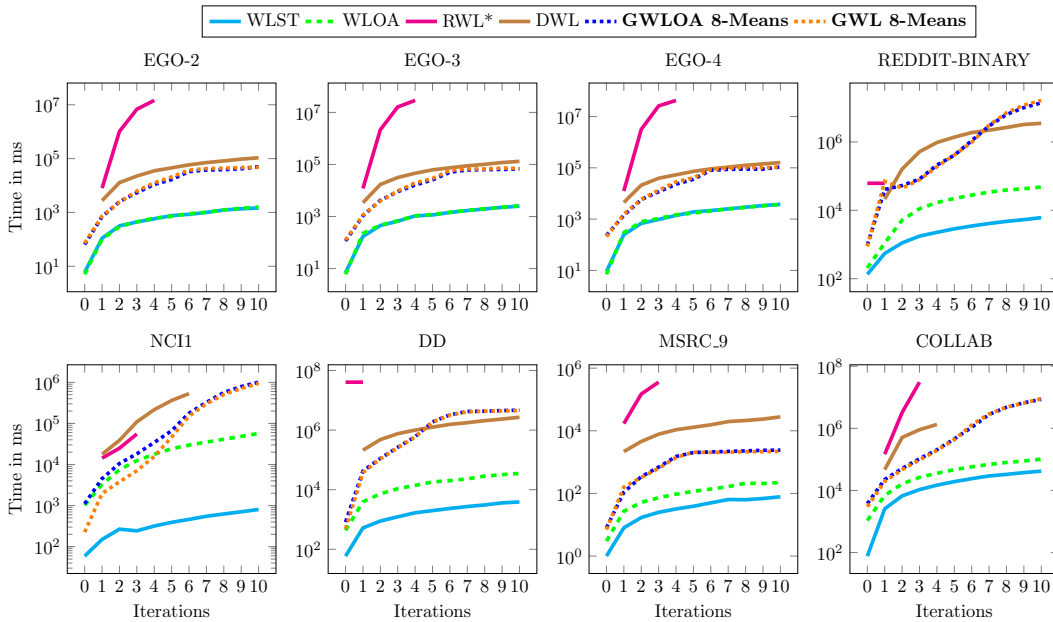


Figure 8: Running time in milliseconds for computing the feature vectors using the different methods. Note that RWL* uses multi-threading, while the other methods do not. Missing values for RWL* and DWL in the larger datasets are due to timeout.

480 **E Running Time - Further Results**

481 Figure 8 shows the running time results for the remaining datasets. While the runtime of our approach
 482 exceeds that of DWL on some of the larger datasets, it enhances the classification accuracy a lot.

483 **F Results on Synthetic Datasets**

484 Table 5 shows the classification accuracy of the different methods on the synthetic datasets (generated
 485 as described in Appendix B), with the best accuracy for each dataset being marked in bold. We can
 486 see, while all kernels can perfectly learn on the datasets without noise, neither WLST, WLOA nor
 487 DWL can manage the noise included in the other datasets, having worse accuracy with increasing
 488 noise. While the decrease in accuracy with decreasing the edge probability is slightly worse than that
 489 of RWL*, our approach has a much lower running time.

Table 5: Average classification accuracy and standard deviation on the synthetic datasets.

Kernel	S_1_0	S_1_10	S_1_20	S_1_50	S_1_100	S_0.8_0	S_0.6_0	S_0.4_0	S_0.2_0
WLST	100.00±0.00	98.68 ±0.34	61.93 ±1.08	54.55 ±0.68	49.78 ±1.18	50.65 ±1.57	48.10 ±1.10	51.98 ±1.40	42.65 ±1.80
DWL	100.00±0.00	98.70 ±0.31	62.10 ±0.98	43.83 ±1.66	51.40 ±0.94	49.80 ±2.01	46.88 ±1.89	49.05 ±1.98	42.85 ±1.98
RWL*	100.00±0.00	100.00±0.00	100.00±0.00	out of time	100.00±0.00	100.00±0.00	99.35 ±0.17	81.93 ±1.00	56.33 ±2.48
WLOA	100.00±0.00	97.65 ±0.44	60.85 ±1.65	47.50 ±1.83	50.23 ±1.24	49.45 ±1.47	48.08 ±1.92	43.23 ±1.21	50.53 ±1.92
GWL	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	92.20 ±0.97	72.53 ±1.78	53.58 ±1.71
GWLOA	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	94.95 ±0.59	72.03 ±1.66	50.90 ±2.63

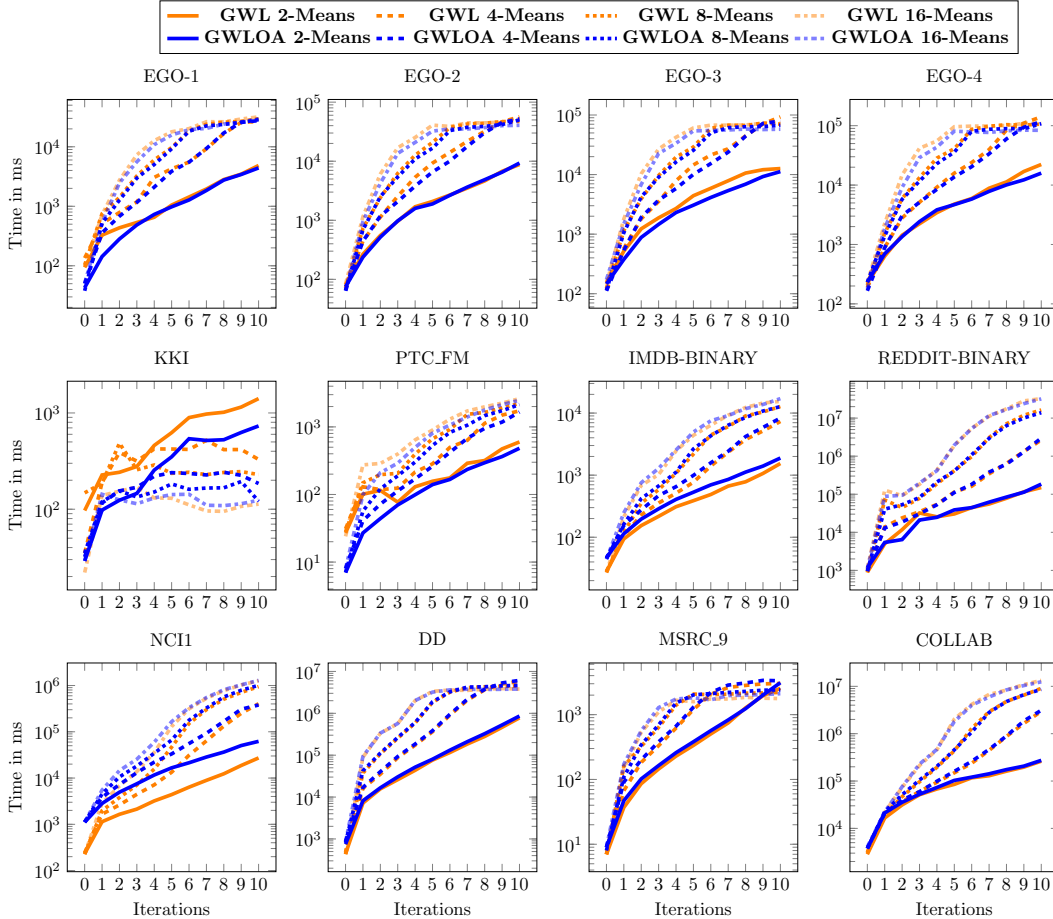


Figure 9: Running time in milliseconds for computing the feature vectors using the different values for parameter k on our newly proposed methods.

490 G Influence of Parameter k on Running Time

491 We investigate which effect the choice of k has on the running time. Figure 9 shows the time needed
 492 for computing the feature vectors using our kernels with $k \in \{2, 4, 8, 16\}$. The difference in running
 493 time between GWL and GWLOA is only marginal on most datasets, only on *KKI* and *NCI1* a larger
 494 difference can be seen. As expected, the running time of both kernels increases with increasing k .
 495 Interestingly, for larger k , the running time does not increase much anymore, after a certain number
 496 of iterations, this might be because the stable coloring was reached by then.

497 H Results on Larger Datasets

498 GWL allows to generate explicit sparse feature vectors just as WLST. Therefore, these kernels can be
 499 used with a linear SVM, which is more efficient than a kernel SVM and makes the application to
 500 larger datasets feasible. We performed additional experiments with these kernels on larger synthetic

Table 6: Parameters used for generating the larger datasets.

Dataset	p	m	Graphs	b	r
$L1$	1	200	10000	25	10
$L2$	1	100	20000	25	10
$L3$	1	200	20000	10	15
$L4$	1	400	20000	25	10

Table 7: Results on larger datasets with running time in seconds and $|f|$ being the number of features (for WLST we only give the order of magnitude).

Kernel	$L1$			$L2$			$L3$			$L4$		
	time	acc	$ f $	time	acc	$ f $	time	acc	$ f $	time	acc	$ f $
WLST	25	50.92	10^6	59	84.57	10^6	45	49.91	10^6	124	49.65	10^6
GWL	208	99.98	61	237	100.0	15	103	100.0	7	137	99.98	7

501 datasets using the linear SVM implementation LIBLINEAR [26]. The datasets were generated using
 502 the same method as described in Appendix B, but with different values for the number of vertices
 503 in the seed graphs, b , and number of vertices, each vertex in the seed graph is replaced by, r (see
 504 Table 6 for the values of the parameters). The experimental setup used in these experiments was as
 505 follows: We split the dataset randomly into a training, validation and test set and chose the parameter
 506 $C \in \{0.1, 1, 10\}$. We used $h \in \{1, \dots, 5\}$ for both kernels and set $k = 2$ for all datasets.

507 In Table 7 we report the running time and the number of features f obtained for the choice of h that
 508 gave the best classification accuracy (for GWL this choice was $h = 2$ for $L4$ and $L5$, $h = 3$ for $L2$
 509 and $h = 5$ for $L1$, for WLST it was $h = 2$ for all datasets except $L4$, on which it was $h = 5$). We
 510 noticed that running the SVM with the feature vectors generated by WLST took much more time
 511 than GWL, which can be explained by the number of features each algorithm produced: For GWL
 512 the number is restricted by the choice of parameters, but for WLST the number of features is only
 513 restricted by the number of vertices in the dataset. In WLST the number of features exceeded 10^6 for
 514 $h \geq 2$. For a small increase in running time in generating the feature vectors, GWL provides not only
 515 a much better classification accuracy than WLST, it also generates less features, but more meaningful
 516 ones.

517 I Approximating the Graph Edit Distance using Optimal Assignments

518 The graph edit distance (GED), a commonly used distance measure for graphs, is defined as the
 519 cost of transforming one graph into the other using edit operations, i.e. deleting or inserting an
 520 isolated vertex or an edge, or relabeling any of the two. An edit path between G and H is a sequence
 521 (e_1, e_2, \dots, e_k) of edit operations that transforms G into H . This means, that applying all operations
 522 in the edit path to G , yields a graph G' that is isomorphic to H .

Edit operations have non-negative costs assigned to them by a cost function c and the GED of two graphs is defined as the cost of a cheapest edit path between them:

$$\text{GED}(G, H) = \min \left\{ \sum_{i=1}^k c(e_i) \mid (e_1, \dots, e_k) \in \Upsilon(G, H) \right\},$$

523 where $\Upsilon(G, H)$ denotes the set of all possible edit paths from G to H .

524 Since computing the GED is NP-hard [27], it is often approximated; often an optimal assignment
 525 between the vertices of the graphs is computed and a (suboptimal) edit path is derived from this
 526 assignment [23]. Figure 10 shows two graphs, an assignment between their vertices indicated by
 527 matching numbers, and a corresponding edit path.

528 If the cost function is a tree metric, such an assignment can be computed in linear time [22] (as
 529 opposed to cubic runtime for arbitrary cost functions). The color hierarchy produced by for example
 530 GWLT can be interpreted as such a tree metric, which means it can be used to find an optimal

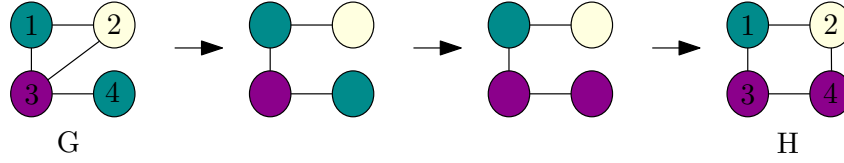


Figure 10: Two graphs G and H with an assignment between their vertices (vertices with matching numbers are assigned to each other) and an edit path derived from this assignment.

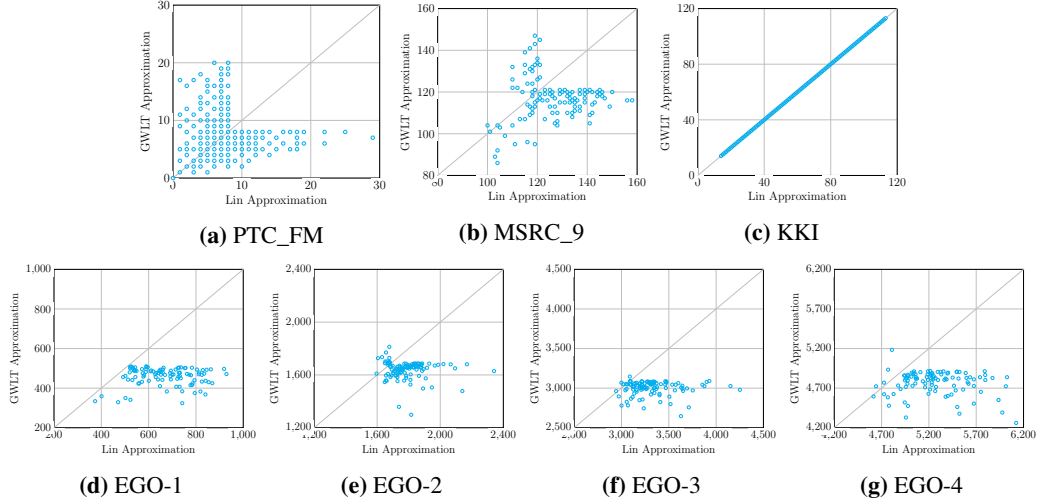


Figure 11: Approximation quality of GWLT compared to Lin [22]. Marks below the gray diagonal indicate that GWLT had a better approximation, while marks above indicate the same for Lin.

531 assignment between the vertices and in turn approximate the graph edit distance. Given a tree T ,
 532 representing a tree metric, an optimal assignment between two sets A and B can be computed as
 533 follows: Let ϕ be a function that associates all elements with their node in the tree (for GWLT this
 534 means, that we associate each vertex with its "newest" color). Then deconstruct T fully (until there
 535 are no leaves left) by repeating these steps:

- 536 1. Pick a random leaf l .
- 537 2. Assign as many elements of A associated with l to elements of B associated with l as possible.
- 538 3. Re-associate remaining elements to the parent of l .
- 539 4. Delete l from T .

540 The resulting assignment is optimal and can be used to derive an edit path as above.

541 J Comparison of Approximation Quality

542 We investigate the accuracy of our method further, by comparing the approximation of Lin [22] to
 543 our method GWLT directly. Figure 11 shows the approximations in comparison. Only graph pairs,
 544 on which at least one of the methods returned a distance below a cutoff threshold are depicted to lay
 545 an emphasis on the more important cases of graphs that are similar.