

A Theoretical and Empirical Investigation into the Equivalence of GNNs and 1-WL+NN

From the faculty of Mathematics, Physics, and Computer Science for the purpose of obtaining the
academic degree of Bachelor of Sciences.

Eric Tillmann Bill

Supervision:

Prof. Dr. rer. nat. Christopher Morris

Informatik 6
RWTH Aachen University

Contents

1	Introduction	3
2	Related Work	4
2.1	Graph Neural Networks	4
2.2	Weisfeiler and Leman Algorithm	4
2.3	Connections between GNNs and the WL algorithm	5
3	Preliminaries	5
3.1	Graph Framework	5
3.2	Permutation-invariance and -equivariance	6
3.3	Weisfeiler and Leman Algorithm	6
3.4	1-WL+NN Framework	7
3.5	Question	7
3.6	1-WL+NN Framework	8
3.7	Graph Neural Networks (Message Passing)	8
4	Main Part	9
4.1	Part I: Theoretical Proof of the Equivalence	9
4.2	Part II: Empirical Analysis and Comparison	9

1 Introduction

Graphs are ubiquitous in various fields of life. Despite not always being explicitly identified as such, the graph data model’s flexibility and simplicity make it an effective tool for modeling a diverse range of data. This includes unexpected instances, such as modeling text or images as a graph, as well as more complex instances like chemical molecules, citation networks, or connectivity encodings of the world wide web (??). Although machine learning has achieved remarkable success with images and text in the last decade, the lack of a significant breakthrough in machine learning for graphs can be attributed to the model’s inherent flexibility and simplicity. While nodes in simple applications may be organized sequentially, as in text, or in a grid-like fashion, as in images, information in graphs can be arranged in more complex ways such as trees, acyclic graphs, or cyclic graphs. This complexity presents a significant challenge in developing a general machine-learning framework capable of accommodating all forms of graph inputs.

In recent years, there has been a significant amount of research conducted to investigate Graph Neural Networks (short GNNs). Among the most promising approaches are those utilizing the message-passing architecture, which was introduced by ? and ?. Empirically, this framework has demonstrated strong performance across various applications (???). However, its expressiveness is limited, as has been proved by the works of ?, as well as ?. These works establish a connection to the commonly used Weisfeiler-Leman¹ algorithm (short 1-WL), originally proposed by ? as a simple heuristic for the graph isomorphism problem. In particular, it has been proven that a GNN based on the message-passing architecture can at most be as good as the 1-WL algorithm in distinguishing graphs. Furthermore, the 1-WL method demonstrates numerous similarities with the fundamental workings of the GNN architecture. It is therefore commonly believed that both methods are to some extent equivalent in their capacity.

In this work, we introduce a novel framework, which we coined ”1-WL+NN,” which involves applying the 1-WL algorithm to an input graph, followed by further processing of the resulting information using a feed-forward neural network. Thereby, we obtain a trainable framework that is suited for all kinds of graph-related tasks, such as graph classification, node regression, and more.

In this thesis, we will conduct an extensive analysis of the newly proposed framework ”1-WL+NN” and its connection to GNNs. The work will be divided into two main parts. Firstly, we will establish the equivalence of both frameworks in terms of their expressiveness. Specifically, we want to show that for every function computed by a GNN, there exists a suitable model based on the ”1-WL+NN” framework with optimal parameters that compute the same function.

Secondly, we will undertake an empirical analysis of the performance of our proposed framework, by testing various configurations and comparing them to state-of-the-art GNN models, that are based on the message-passing architecture. Additionally, we will explain potential observable differences between the two frameworks, despite their shared similarities. Finally, the thesis will conclude with recommendations for further research based on these findings.

¹Based on <https://www.iti.zcu.cz/wl2018/pdf/leman.pdf>, we will use the spelling ”Leman” here, as requested by A. Leman, the co-inventor of the algorithm.

2 Related Work

2.1 Graph Neural Networks

In recent years, machine learning models have experienced a surge in popularity due to their significant performance advantages over conventional methods and their ability to autonomously adapt to their tasks. However, a closer examination of the applications where these models are been used reveals that they are highly specialized for the specific input of each application. For instance, modern convolutional neural networks (short CNNs) are designed to take in fixed-sized, grid-like data structures such as images, while modern language models process sequential data like textfiles.

The relevance of graphs to these examples lies in the fact that graphs can be used to model various types of inputs across many applications, and they provide a more general framework for modeling data. To illustrate, an image can be modeled as a graph for a CNN, where each pixel corresponds to a node in the graph holding the brightness value for each color value, and each node is connected to its neighboring pixels. Similarly, for sequential data like textfiles, one can encode a directed graph where each word in this file is represented as a node with the word as a feature, and it is connected outgoingly to the next following word. With these examples, we wanted to highlight the flexibility of how graphs can model data, however, this is also problematic, as this makes it particularly hard to construct a general machine-learning model on graphs. Levering any constraints on the format or size of the input can significantly limit the model’s generality, and since graphs sizes and formats can vary within their applications, e.g. classification of molecules (PUBCHEM ?), the need for a general model is of great interest.

From the work of ?, as well ?, the so-called message-passing architecture emerged for Graph Neural Networks (short GNNs). This can be understood as a framework that never changes its input graph structurally and only modifies the node’s features in each layer. In more detail, in each layer, a GNN based on the message-passing architecture, computes for each node a new feature, based on its current feature and the features of its neighbors. Later in section 3.7, we will give a more formal definition of this architecture. Throughout this thesis, I will use the term GNN and message-passing architecture interchangeably.

2.2 Weisfeiler and Leman Algorithm

The (1-dimensional) Weisfeiler-Leman algorithm (short 1-WL), proposed by ?, was initially designed as a simple heuristic for the *graph isomorphism problem*, but due to its interesting properties, its simplicity, and its good performance, the 1-WL algorithm gained a lot of attention from researchers across many fields. One of the most noticeable of these properties is, that the algorithm color codes the nodes of the input graph in such a way, that in each iteration, each color encodes a local learned substructure.

It works by coloring all nodes in each iteration the same color that fulfill two properties: 1. the nodes already share the same color and 2. the frequencies of the colors of their neighbors are equal. The algorithm continues as long as the number of colors changes in each iteration. For determining whether two graphs are non-isomorphic, the heuristic applies the algorithm to both graphs simultaneously and concludes that the graphs are non-isomorphic as soon as the number of occurrences of a color is different between the graphs. We present a more formal definition of the algorithm later in the section 3.3.

Since the *graph isomorphism problem* is difficult to solve due to it being NP-complete, the 1-WL algorithm, running in polynomial deterministic time, cannot solve the problem completely.

Moreover, [?] have constructed counterexamples of non-isomorphic graphs that the heuristic fails to distinguish, e.g. figure 1. However, following the work of [?], this simple heuristic is still quite powerful and has a very low probability of failing to distinguish non-isomorphic graphs when both graphs are uniformly chosen at random.

To overcome the limited expressiveness of the 1-WL algorithm, it was generalized to the k -dimensional Weisfeiler-Leman algorithm (short k -WL), which works with k -tuples over the set of nodes. Interestingly, this created a hierarchy for the expressiveness of determining non-isomorphism, such that for all $k \in \mathbb{N}$ there exists a pair of non-isomorphic graphs that can be distinguished by the $(k + 1)$ -WL but not by the k -WL ([?]).

2.3 Connections between GNNs and the WL algorithm

A connection between GNNs based on the message-passing architecture and the 1-WL algorithm seems quite natural since both share similar properties in terms of how they process graph data. Most noticeably, both methods never change the graph structurally, since they only compute new node features in each iteration. And additionally, both methods use a 1-hop neighborhood aggregation as their basis for the computation of the new node feature. Following this intuition of both methods being very similar, many authors showed a theoretical connection between these methods. [?], as well as [?], showed that GNN's expressiveness power is upper bounded by the 1-WL in terms of distinguishing non-isomorphic graphs. In addition, [?] also proposed a new k -GNN architecture that works over the set of subgraphs of size k . Interestingly, the authors showed that the proposed hierarchy over $k \in \mathbb{N}$ is equivalent to the k -WL hierarchy in terms of their expressive in distinguishing non-isomorphic graphs, meaning if there exists a k -GNN that can distinguish two non-isomorphic graphs then it is equivalent to say that the k -WL algorithm can distinguish these graphs as well.

3 Preliminaries

We first introduce a couple of notations that will be used in this thesis. With $[n]$, we denote the set $\{1, \dots, n\} \subset \mathbb{N}$ for any $n \in \mathbb{N}$ and with $\{\!\!\{ \dots \}\!\!\}$ we denote a multiset which is formally defined as a 2 tuple (X, m) with X being a set of all unique elements and $m : X \rightarrow \mathbb{N}_{\geq 1}$ a mapping that maps every element in X to its number of occurrences in the multiset.

3.1 Graph Framework

A graph is denoted by G and is a 3 tuple $G := (V, E, l)$ that consists of the set of all nodes V , the set of all edges $E \subseteq V \times V$ and a label function $l : M \rightarrow \Sigma$ with M being either $V, V \cup E$ or E and $\Sigma \subset \mathbb{N}_0$ a finite alphabet. Moreover, let \mathcal{G} be the set of all graphs. Note, that our definition of the label function allows for graphs with labels either only on the nodes, only on the edges, or on both nodes and edges. Sometimes the values assigned by l are called features, but this is usually only the case when Σ is multidimensional, which we do not cover in this thesis. In addition, although we have defined it this way, the labeling function is optional, and in cases where no labeling function is given, we add the trivial labeling function $f_0 : V(G) \rightarrow \{0\}$. Further, G can be either directed or undirected, depending on the definition of E , where $E \subseteq \{(v, u) \mid v, u \in V\}$ defines a directed and $E \subseteq \{(v, u), (u, v) \mid v, u \in V, v \neq u\}$ defines an undirected graph. Additionally, we will use the notation $V(G)$ and $E(G)$ to denote the set of

nodes of G and the set of edges of G respectively. With $\mathcal{N}(v)$ for $v \in V(G)$ we denote the set of neighbors of v with $\mathcal{N}(v) := \{u \mid (u, v) \in E(G)\}$.

A coloring of a Graph G is a function $C : V(G) \rightarrow \mathbb{N}_0$ that assigns each node in the graph a color (here a positive integer). Further, a coloring C induces a partition P on the set of nodes, for which we define C^{-1} being the function that maps each color $c \in \mathbb{N}_0$ to its class of nodes with $C^{-1}(c) = \{v \in V(G) \mid C(v) = c\}$. In addition, we let $h_{G,C} = \{\{C(v) \mid v \in V(G)\}\}$ be the histogram of graph G with coloring C , that contains for every color in the image of $V(G)$ under C the color and its frequency.

3.2 Permutation-invariance and -equivariance

We use S_n to denote the symmetric group over the elements $[n]$ for any $n > 0$. S_n consists of all permutations over these elements. Let G be a graph with $V(G) = [n]$, applying a permutation $\pi \in S_n$ on G , is defined as $G_\pi := \pi \cdot G$ where $V(G_\pi) = \{\pi(1), \dots, \pi(n)\}$ and $E(G_\pi) = \{(\pi(v), \pi(u)) \mid (v, u) \in E(G)\}$. We will now introduce two key concepts for classifying functions on graphs. Let $f : \mathcal{G} \rightarrow \mathcal{X}$ be an arbitrary function and let $V(G) = [n]$ for some $n \in \mathbb{N}$:

1. The function f is *permutation-invariant* if and only if for all $G \in \mathcal{G}$ where $n_G := |V(G)|$ and for every $\pi \in S_{n_G}$: $f(G) = f(\pi \cdot G)$.
2. The function f is *permutation-equivariant* if and only if for all $G \in \mathcal{G}$ where $n_G := |V(G)|$ and for every $\pi \in S_{n_G}$: $f(G) = \pi^{-1} \cdot f(\pi \cdot G)$

3.3 Weisfeiler and Leman Algorithm

The Weisfeiler-Leman algorithm consists of two main parts, first the coloring algorithm and second the graph isomorphism test. We will introduce them in this section in this order.

The Weisfeiler-Leman graph coloring algorithm

Let $G = (V, E, l)$ be a graph, then in each iteration i , the 1-WL computes a node coloring $C_i : V(G) \rightarrow \mathbb{N}$, which depends on the coloring of the neighbors and the node itself. In iteration $i = 0$, the initial coloring is $C_0 = l$ or if l is non existing $C_0 = c$ for an arbitrary constant $c \in \mathbb{N}$. For $i > 0$, the algorithm assigns a color to $v \in V(G)$ as follows:

$$C_i(v) = \text{RELABEL}((C_{i-1}(v), \{\{C_{i-1}(u) \mid u \in \mathcal{N}(v)\}\}))$$

Where RELABEL injectively maps the above pair to a unique, previously not used, natural number. The algorithm terminates when the number of colors between two iterations does not change, meaning the algorithm terminates after iteration i if the following condition is satisfied:

$$\forall v, w \in V(G) : C_i(v) = C_i(w) \iff C_{i+1}(v) = C_{i+1}(w) \quad (0.1)$$

Upon terminating we define $C_\infty := C_i$ as the stable coloring. The algorithm always terminates after $n_G := |V(G)|$ iterations (?). Moreover, based on the work of ? about efficient refinement strategies, ? proved that the stable coloring C_∞ can be computed in time $\mathcal{O}(|V(G)| + |E(G)| \cdot \log |V(G)|)$.

The Weisfeiler-Leman Graph Isomorphism Test

To determine if two graphs $G, H \in \mathcal{G}$ are non-isomorphic (short $G \not\cong H$), one applies the 1-WL coloring algorithm on both graphs "in parallel" and checks after each iteration if the occurrences of each color are equal, else the algorithm would terminate and conclude non-isomorphic. Formally, the algorithm concludes non-isomorphic in iteration i if there exists a color c such that:

$$|\{v \in V(G) \mid c = C_i(v)\}| \neq |\{v \in V(H) \mid c = C_i(v)\}| \quad (0.2)$$

Note that this test is only sound and not complete for the problem of graph isomorphism. Counterexamples where the algorithm fails to distinguish non-isomorphic graphs can be easily constructed, see Figure 1 which was discovered and proven by ?.

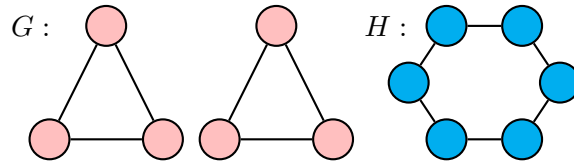


Figure 1: An example of two graphs G and H that are non-isomorphic but cannot be distinguished by the 1-WL

3.4 1-WL+NN Framework

Let $\mathcal{X} \subset \mathcal{G}$ be the set of all graphs with n nodes. Let F_{1-WL} be the function that works over \mathcal{X} and maps every graph G to its color histogram of the stable coloring computed when applying the 1-WL algorithm. More formally, $F_{1-WL} : \mathcal{X} \rightarrow \mathbb{N}^c$

3.5 Question

Let \mathcal{X} be the set of all labeled graphs G with n nodes.

$F : \mathcal{X} \rightarrow R^C$ the function that maps every graph $G \in \mathcal{X}$ to its 1-WL color histogram.

The class 1-WL+NN of functions consists of all functions f with the format: $f : \text{MLP} \circ F$

Wie kann man den das Color-Histogramm in einen Vektor aus R^C mappen?

Meine Idee dazu wäre, $R = \mathbb{N}$ und für jede Farbe $i \in \mathbb{N}$ setzten wir die i .te Komponente des Vektors auf die Anzahl dieser Farbe im Histogramm. Zum Beispiel wenn wir das Color-Histogramm $h = \{1, 2, 2, 4\}$ haben, dann würden wir auf den Vektor: $(1, 2, 0, 1, 0, \dots, 0)^T \in \mathbb{N}^C$ abbilden.

In diesem Verfahren liese sich C ganz gut abschätzen. Problematisch wird es nur wenn wir labeled graphen $G = (V, E, l)$ betrachten, vorallem wenn $\forall v \in V(G) : l(v) > C$ und jedes label einzigartig ist (d.h. $|V(G)| = |l(V(G))|$), dann würde der 1-WL algorithmus nach der 1. Iteration bereits terminieren und als Color-Histogramm $\{l(v) \mid v \in V(G)\}$ liefern. Eine Beschränkung der Labels auf $\{1, \dots, C\}$, oder ein Mapping auf neue Labels zwischen $\{1, \dots, C\}$ ist natürlich möglich, wäre aber leider schlecht für späteres Testen mit echten Testcases, da dann ein Verlust des Informationsgehalts droht. Oder man liese den 1-WL Alogrithmus eine weitere Iteration laufen, sodass dieser ein stabiles Coloring liefert dessen Bild in $\{1, \dots, C\}$ liegt. Wie soll dieser Fall behandelt werden, oder ist dies garkein Problem?

3.6 1-WL+NN Framework

Let \mathcal{Y} be the task-specific output set (e.g. set of class labels), \mathcal{NN} a feedforward neural network, enc an encoding function, π_{pool} a pooling function. Further, for $G \in \mathcal{G}$, let $(C_\infty^i)_G$ be the final coloring upon termination when applying the 1-WL algorithm on G . Then the computed function \mathcal{A} is:

$$\mathcal{A} : \mathcal{X} \rightarrow \Sigma, G \mapsto \mathcal{NN} \circ \pi_{\text{pool}}(\{(C_\infty^i)_G(v) \mid v \in V(G)\}) \quad (0.3)$$

3.7 Graph Neural Networks (Message Passing)

Let $G = (V, E, l)$ be an arbitrary graph. A Graph Neural Network (GNN) is a composition of multiple layers where each layer t passes a vector representation of each node v or edge e through $f^{(t)}(v)$ or $f^{(t)}(e)$ respectively and retrieves thereby a new graph that is structurally identical but has new feature information. Note that in the following we will restrict the definition to only consider node features, however, one can easily extend it to also include edge features.

To begin with, we need a function $f^{(0)} : V(G) \rightarrow \mathbb{R}^{1 \times d}$ that is consistent with l , that translates all labels into a vector representation. Further, for every $t > 0$, f is of the format:

$$f^{(t)}(v) = f_{\text{merge}}^{W_{1,t}}(f^{(t-1)}(v), f_{\text{agg}}^{W_{2,t}}(\{f^{(t-1)}(w) \mid w \in \mathcal{N}(v)\})) \quad (0.4)$$

Where $f_{\text{merge}}^{W_{1,t}}$ and $f_{\text{agg}}^{W_{2,t}}$ are arbitrary differentiable functions with $W_{1,t}$ and $W_{2,t}$ their respective parameters. Additionally, $f_{\text{agg}}^{W_{2,t}}$ has to be permutation-invariant. To demonstrate what kind of functions are typically used, we provide the functions used by ?:

$$f_{\text{merge}}^{W_{1,t}}(v, \text{agg}) = \sigma(W^t \times \text{concat}(v, \text{agg})) \quad (0.5)$$

$$f_{\text{agg}}^{W_{2,t}} = \max(\{\sigma(W_{\text{pool}}^t \times f^{(t-1)}(u) + b \mid u \in \mathcal{N}(v)\})\}) \quad (0.6)$$

Where σ is a non-linear elementwise activation function; W^k, W_{pool} and b are trainable parameters and concat the concatenation function.

Depending on the objective, whether the GNN is tasked with a graph or only a node or edge task, the last layer differs. In the case of graph tasks, we add a permutation-invariant aggregation function to the end, here called **READOUT**, that aggregates over every node and computes a fixed-size output vector for the entire graph, e.g. a label for graph classification. In order to ensure that we can train the GNN in an end-to-end fashion, we require **READOUT** to be also differentiable.

Let \mathcal{A} be an instance of the described GNN framework. Further, let $K \in \mathbb{N}$ be the number of layers of the GNN, \mathcal{G} the set of all graphs, \mathcal{Y} the task-specific output set (e.g. labels of a classification task), then the overall function computed by \mathcal{A} is:

$$\mathcal{A} : \mathcal{G} \rightarrow \mathcal{Y} : x \mapsto f^{(K)} \circ \dots \circ f^{(0)}(x) \quad (0.7)$$

$$\mathcal{A} : \mathcal{G} \rightarrow \mathcal{Y} : x \mapsto \text{READOUT} \circ f^{(K)} \circ \dots \circ f^{(0)}(x) \quad (0.8)$$

As we require all aggregation functions to be permutation-invariant, the total composition \mathcal{A} is permutation-invariant, and similarly, it is also differentiable. This enables us to train \mathcal{A} like any other machine learning method in an end-to-end fashion, regardless of the underlying encoding used for graphs. This definition and use of notation are inspired by ? and ?.

4 Main Part

In this section, we will discuss the topic of this thesis. The thesis will be about the connection between graph neural networks based on the message-passing architecture and the 1-dimensional Weisfeiler-Leman algorithm. In particular, we will show a theoretical equivalence between GNNs and 1-WL+NN about their expressive in the first part. The second part will then be based on the results of the first part, an empirical analysis of the performance of GNNs and 1-WL+NN in different configurations. We will now introduce each part individually.

4.1 Part I: Theoretical Proof of the Equivalence

In this part, we will start by introducing the framework we coined 1-WL+NN, and demonstrate how a model of this framework can be used and trained in an end-to-end fashion for graph tasks. More importantly, we will build a connection to GNNs by trying to prove the following hypothesis:

For every function \mathcal{A} computed by a GNN (definition in section 3.7), there exists a 1-WL+NN model (definition in section 3.4) computing \mathcal{A} as well.

After proving this hypothesis formally, we can conclude that GNNs and 1-WL+NN models share the same capacity.

As of today, no known research to us investigates the relationship between these two frameworks. Nonetheless, several works, as outlined in 2.3, offer results that hint that this hypothesis may be valid. Additionally, in the research conducted by ?, the author empirically demonstrated how well the 1-WL test is in distinguishing non-isomorphic graphs across 9 datasets and used these results to give an upper bound on the actual classification task of the datasets when training an individual GNN on each dataset. This demonstrates, how the expressiveness of a graph algorithm on an arbitrary task is somehow limited by its capacity in distinguishing non-isomorphism.

Our approach for showing the validness of the hypothesis is, that the proving direction of "1-WL+NN \subseteq GNN" is relatively trivial, as can easily encode the 1-WL coloring algorithm in each GNN layer. The other direction, however, showing "GNN \subseteq 1-WL+NN", is more challenging, for which we want to take inspiration from the proof presented in section 3.3 by ?.

4.2 Part II: Empirical Analysis and Comparison

In this part, we will provide a comprehensive analysis of the performance of the 1-WL+NN framework by testing it on well-established benchmark datasets for GNNs. With this analysis we will try to answer the following questions:

- Q1) Which encoding of the feature space for 1-WL+NN framework has the best performance in generalizing? Does this result align with the results of research in other fields?
- Q2) Is there a performance difference between both frameworks, 1-WL+NN and GNNs, in generalizing? And if so, which one generalizes better after fewer training iterations? Is there an explanation for this behavior?
- Q3) Is one of the tasks better suited for a specific task setting? For example, is more suitable for graph classification? Why could this be, what is the fundamental difference between both frameworks leading to this result?

As for the model configurations, we want to test, we decided on the following 6 configurations, where the first three are 1-WL+NN models, and the last three are state-of-the-art GNNs:

1. 1-WL+NN model using a one-hot encoding of its feature space
2. 1-WL+NN model using a look-up-table for the encoding of its feature space
3. 1-WL+NN model using a GNN for the encoding of its feature space
4. Graph Convolutional Network (short GCN) by ? with and without training
5. GraphSAGE developed by ? with and without training
6. Graph Isomorphism Network (short GIN) developed by ? with and without training

For each configuration, we will choose suitable parameters, as well as leaving us the option of using maybe other configurations or GNNs models. Since this analysis completely depends on the results of the first part of this thesis, such that we can only speculate now, about suitable test cases.

For running the test cases, we will implement the 1-WL+NN with all its different configurations in Python using the open source library PYTORCH² and the open source extension PYTORCH GEOMETRIC³.

We will select the datasets to be used for benchmarking from the TU-DATASET, a curated collection of graph datasets that are highly suitable for training graph-based algorithms. This collection offers data from diverse applications of varying sizes, enabling us to thoroughly evaluate the performance of our frameworks on a wide range of inputs. This dataset is the result of extensive work by ?.

²Open source machine learning framework that was originally developed by Meta AI and does now belong to the Linux Foundation umbrella. <https://pytorch.org>

³Open source library that acts as an extension to PyTorch and allows for easy writing and training of graph neural networks. <https://pytorch-geometric.readthedocs.io>