



HAL
open science

Graph Ordering Attention Networks

Michail Chatzianastasis, Johannes Lutzeyer, George Dasoulas, Michalis Vazirgiannis

► **To cite this version:**

Michail Chatzianastasis, Johannes Lutzeyer, George Dasoulas, Michalis Vazirgiannis. Graph Ordering Attention Networks. AAAI Conference on Artificial Intelligence (AAAI), Feb 2023, Washington DC, United States. pp.7006-7014, 10.1609/aaai.v37i6.25856 . hal-04447535

HAL Id: hal-04447535

<https://hal.science/hal-04447535>

Submitted on 8 Feb 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Graph Ordering Attention Networks

Michail Chatzianastasis^{1*}, Johannes F. Lutzeyer¹, George Dasoulas¹ and Michalis Vazirgiannis¹

¹ DaSciM, LIX, École Polytechnique, Institut Polytechnique de Paris, France.

{michail.chatzianastasis, johannes.lutzeyer}@polytechnique.edu,
george.dasoulas1@gmail.com, mvazirg@lix.polytechnique.fr

Abstract

Graph Neural Networks (GNNs) have been successfully used in many problems involving graph-structured data, achieving state-of-the-art performance. GNNs typically employ a message-passing scheme, in which every node aggregates information from its neighbors using a permutation-invariant aggregation function. Standard well-examined choices such as the mean or sum aggregation functions have limited capabilities, as they are not able to capture interactions among neighbors. In this work, we formalize these interactions using an information-theoretic framework that notably includes synergistic information. Driven by this definition, we introduce the Graph Ordering Attention (GOAT) layer, a novel GNN component that captures interactions between nodes in a neighborhood. This is achieved by learning local node orderings via an attention mechanism and processing the ordered representations using a recurrent neural network aggregator. This design allows us to make use of a permutation-sensitive aggregator while maintaining the permutation-equivariance of the proposed GOAT layer. The GOAT model demonstrates its increased performance in modeling graph metrics that capture complex information, such as the betweenness centrality and the effective size of a node. In practical use-cases, its superior modeling capability is confirmed through its success in several real-world node classification benchmarks.

1 Introduction

Graph Neural Networks (GNNs) achieve remarkable success in machine learning problems on graphs [Scarselli et al., 2009, Kipf and Welling, 2017, Bronstein et al., 2021]. In these problems, data arises in the structure of attributed graphs, where in addition to the node and edge sets defining a graph, a set of feature vectors containing data on each node is present. The majority of GNNs learn node representations using a message-passing scheme [Gilmer et al., 2017].

In such message passing neural networks (MPNN) each node iteratively aggregates the feature vectors or hidden representations of its neighbors to update its own hidden representation. Since no specific node ordering exists, the aggregator has to be a permutation-invariant function [Xu et al., 2019].

Although MPNNs have achieved great results, they have severe limitations. Their permutation-invariant aggregators treat neighboring nodes as a set and process them individually, omitting potential interactions between the large number of subsets that the neighboring nodes can form. Therefore, current MPNNs cannot observe the entire structure of neighborhoods in a graph [Pei et al., 2020] and cannot capture all *synergistic interactions* between neighbors [Murphy et al., 2019, Wagstaff et al., 2021].

The concept of synergy is important in many scientific fields and is central to our discussion here. It expresses the fact that some source variables are more informative when observed together instead of independently. For example in neuroscience, synergy is observed when the target variable corresponds to a stimulus and the source variables are the responses of different neurons [Bizzi and Cheung, 2013]. Synergistic information is often presented in biological cells, where extra information is provided by patterns of coincident spikes from several neurons [Brenner et al., 2000]. In gene-gene interactions, synergy is present when the contribution of two mutations to the phenotype of a double mutant is larger than the expected additive effects of the individual mutations [Pérez-Pérez et al., 2009]. We believe the consideration of synergistic information to have great potential in the GNN literature.

In this paper, to better understand interactions between nodes, we introduce the Partial Information Decomposition (PID) framework [Williams and Beer, 2010] to the graph learning context. We decompose the information that neighborhood nodes have about the central node into three parts: unique information from each node, redundant information, and synergistic information due to the combined information from nodes. We furthermore show that typical MPNNs cannot capture redundant and synergistic information.

To tackle these limitations we propose the *Graph Ordering Attention (GOAT)* layer, a novel architecture that can capture all sources of information. We employ self-attention to construct a permutation-invariant ordering of the nodes in each neighborhood before we pass these ordered sequences

*Contact Author

to a Recurrent Neural Network (RNN) aggregator. Using a permutation-sensitive aggregator, such as the Long Short-Term Memory (LSTM) model, allows us to obtain larger representational power [Murphy et al., 2019] and to capture the redundant and synergistic information. We further argue that the ordering of neighbors plays a significant role in the final representation [Vinyals et al., 2016] and demonstrate the effectiveness of GOAT versus other non-trainable and/or permutation-sensitive aggregators with a random ordering [Hamilton et al., 2017].

Our main contributions are summarized as follows:

1. We present a novel view of learning on graphs based on information theory and specifically on the Partial Information Decomposition. We further demonstrate that typical GNNs can not effectively capture redundant and synergistic information between nodes.
2. We propose the *Graph Ordering Attention (GOAT)* layer, a novel GNN component that can capture synergistic information between nodes using a recurrent neural network (LSTM) as an aggregator. We highlight that the ordering of the neighbors is crucial for the performance and employ a self-attention mechanism to learn it.
3. We evaluate GOAT in node classification and regression tasks on several real-world and synthetic datasets and outperform an array of state-of-the-art GNNs.

2 Preliminaries and Related Work

We begin by defining our notation and problem context.

Problem Formulation and Basic Notation. Let a graph be denoted by $G = (V, E)$, where $V = \{v_1, \dots, v_N\}$ is the node set and E is the edge set. Let $\mathbf{A} \in \mathbb{R}^{N \times N}$ denote the adjacency matrix, $\mathbf{X} = [x_1, \dots, x_N]^T \in \mathbb{R}^{N \times d_I}$ be the node features and $\mathbf{Y} = [y_1, \dots, y_N]^T \in \mathbb{N}^N$ the label vector. We denote the neighborhood of a vertex u by $\mathcal{N}(u)$ such that $\mathcal{N}(u) = \{v : (v, u) \in E\}$ and the neighborhood features by the multiset $\mathbf{X}_{\mathcal{N}(u)} = \{x_v : v \in \mathcal{N}(u)\}$. We also define the neighborhood of u including u as $\overline{\mathcal{N}}(u) = \mathcal{N}(u) \cup \{u\}$ and the corresponding features as $\mathbf{X}_{\overline{\mathcal{N}}(u)}$. The goal of semi-supervised node classification and regression is to predict the labels of a test set given a training set of nodes.

Graph Neural Networks. GNNs exploit the graph structure \mathbf{A} and the node features \mathbf{X} in order to learn a hidden representation h_u of each node u such that the label y_u can be predicted accurately from h_u [Gori et al., 2005, Scarselli et al., 2009]. Most approaches use a neighborhood message-passing scheme, in which every node updates its representation by aggregating the representations of its neighbors and combining them with its previous representation,

$$m_u^{(l)} = \text{Aggregate}^{(l)} \left(\left\{ h_v^{(l-1)} : v \in \mathcal{N}(u) \right\} \right), \quad (1)$$

$$h_u^{(l)} = \text{Combine}^{(l)} \left(h_u^{(l-1)}, m_u^{(l)} \right),$$

where $h_u^{(l)}$ denotes the hidden representation of node u at the l^{th} layer of the GNN architecture. Note that we often omit the superscript (l) to simplify the notation.

Typically GNNs employ a permutation-invariant ‘‘Aggregate’’ function to yield a permutation-equivariant GNN layer [Bronstein et al., 2021]. Permutation invariance and equivariance will be defined formally now.

Definition 2.1. Let S_M denote the group of all permutations of a set containing M elements. A function $f(\cdot)$ is *permutation-equivariant* if for all $\pi \in S_M$ we have $\pi f(\{x_1, x_2, \dots, x_M\}) = f(\{x_{\pi(1)}, x_{\pi(2)}, \dots, x_{\pi(M)}\})$. A function $f(\cdot)$ is *permutation-invariant* if for all $\pi \in S_M$ we have $f(\{x_1, x_2, \dots, x_M\}) = f(\{x_{\pi(1)}, x_{\pi(2)}, \dots, x_{\pi(M)}\})$.

2.1 Common Aggregators and Their Limitations

We now describe some of the most well-known aggregators and discuss their limitations. Our analysis is based on two important properties that an aggregator should have:

1. **Relational Reasoning:** The label of a node may depend not only on the unique information of each neighbor, but also on the joint appearance and interaction of multiple nodes [Wagstaff et al., 2021]. With the term ‘‘relational reasoning’’ we describe the property of capturing these interactions, i.e., synergistic information, when aggregating neighborhood messages.
2. **Injectivity:** As shown in Xu et al. [2019], a powerful GNN should map two different neighborhoods, i.e., multisets of feature vectors, to different representations. Hence, the aggregator should be injective.

The *mean* and *max* functions are commonly used to aggregate neighborhood information [Kipf and Welling, 2017]. However, they are neither injective [Xu et al., 2019] nor able to perform relational reasoning as they process each node independently. The *summation operator followed by a multilayer perceptron* was recently proposed [Xu et al., 2019]. This aggregator is injective but cannot perform relational reasoning and usually requires a large latent dimension [Wagstaff et al., 2019, 2021]. In the Graph Attention Networks (GAT) [Veličković et al., 2018a], the representation of each node is computed by applying a *weighted summation* of the representations of its neighbors. However, the attention function is not injective since it fails to capture the cardinality of the neighborhood. Recently, an improved version of the GAT was published [Brody et al., 2022] and also, a new type of attention was proposed [Zhang and Xie, 2020], that preserves the cardinality of the neighborhood and therefore is injective. Nevertheless, none of these models can capture interactions between neighbor nodes as each attention score is computed based only on the representations of the central node and one neighbor node. In Section 3.2 we provide further details on why typical aggregators fail, from an information theoretic perspective.

2.2 Permutation-Sensitive Aggregators

Several authors have proposed the use of permutation-sensitive aggregators to tackle the limitations of permutation-invariant aggregators. In particular, Niepert et al. [2016] propose to order nodes in a neighborhood according to some labeling, e.g., the betweenness centrality or PageRank score, to assemble receptive fields, possibly extending beyond the 1-hop neighborhood of a central node, which are then fed to

a Convolutional Neural Network (CNN) architecture. While this approach demonstrates good performance, it relies on the fixed chosen ordering criterion to be of relevance in the context of a given dataset and chosen learning task. Gao et al. [2018] propose to only work with the k largest hidden state values for each hidden state dimension in each neighborhood. While not explicitly ordering the neighboring nodes, this operation summarises any given neighborhood in a fixed size feature matrix and enables the use of permutation-sensitive aggregators, CNNs in their case. Of course the choice of k involves a loss of information in almost all cases, i.e., when k is smaller than the maximal degree in the graph. In the Janosy Pooling [Murphy et al., 2019] approach, a permutation-invariant aggregator is obtained by applying a permutation-sensitive function to all $n!$ permutations. Since the computational cost of this approach is very high, they also propose an approximation, sampling only a limited number of permutations. Similarly, in the GraphSage [Hamilton et al., 2017] model, a random permutation of each neighborhood is considered and then passed to an LSTM. However, it has been observed that even in the graph domain, where typically no natural ordering of nodes is known, there exist some orderings that lead to better model performance [Vinyals et al., 2016]. Whether these high performance orderings are discovered during the training process is left to chance in the GraphSage and Janosy Pooling models. In contrast, our method learns a meaningful ordering of neighbors with low complexity by leveraging the attention scores.

3 An Information Theory Perspective

In this section, we show how neighborhood dependencies can be encoded in the Partial Information Decomposition framework. This decomposition will motivate us to build a more expressive GNN layer, that is able to capture various interactions among neighborhood nodes.

3.1 Partial Information Decomposition

The domain of information theory provides a well-established framework for measuring neighborhood influence. A few graph representation learning results capitalize on information-theoretic tools, either assuming a probability distribution over the feature vectors [Veličković et al., 2018b, Peng et al., 2020] or over the structural characteristics [Luo et al., 2021, Dasoulas et al., 2020].

The majority of GNNs (including the attention-based models) use an aggregation that does not capture interactions among neighbors. Mutual information is a measure that can give us an insight in the omitted informative interactions.

Definition 3.1. For a given node $u \in V$, let $\mathbf{H}_{\bar{\mathcal{N}}(u)} = [h_{v_1}, \dots, h_{v_{|\bar{\mathcal{N}}(u)|}}] \in \mathbb{R}^{|\bar{\mathcal{N}}(u)| \times d}$ denote the hidden representations of the nodes in $\bar{\mathcal{N}}(u)$. Then, if we assume that $\mathbf{H}_{\bar{\mathcal{N}}(u)}$ and h_u follow distributions $p(\mathbf{H}_{\bar{\mathcal{N}}(u)})$ and $p(h_u)$, respectively, the *mutual information* between h_u and $\mathbf{H}_{\bar{\mathcal{N}}(u)}$ is defined as

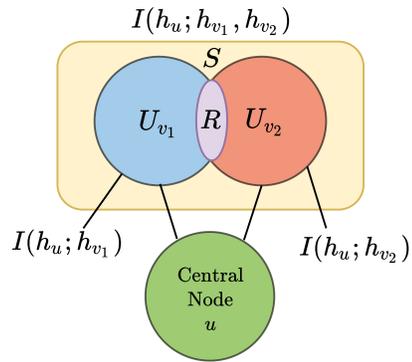


Figure 1: An illustration of the Partial Information Decomposition for the case of one central node u and two neighbors v_1, v_2 . Each of the mutual information terms $I(h_u; h_{v_1})$ and $I(h_u; h_{v_2})$ consists of the unique information provided by v_1 (U_{v_1} , blue patch) and v_2 (U_{v_2} , red patch), respectively, as well as the shared information of v_1 and v_2 (R , purple patch). The joint mutual information $I(h_u; h_{v_1}, h_{v_2})$ (yellow box encompassing the inner two circles) consists of four elements: the unique information in the neighbors v_1 and v_2 , their redundant information and additionally the synergistic information, $I(h_u; h_{v_1}, h_{v_2}) = U_{v_1} + U_{v_2} + R + S$.

$$I(h_u; \mathbf{H}_{\bar{\mathcal{N}}(u)}) = \iint p(h_u, \mathbf{H}_{\bar{\mathcal{N}}(u)}) \log \left(\frac{p(h_u, \mathbf{H}_{\bar{\mathcal{N}}(u)})}{p(h_u)p(\mathbf{H}_{\bar{\mathcal{N}}(u)})} \right) dh_u d\mathbf{H}_{\bar{\mathcal{N}}(u)}. \quad (2)$$

Following Williams and Beer [2010], (2) can be decomposed into three components as follows:

$$I(h_u; \mathbf{H}_{\bar{\mathcal{N}}(u)}) = \sum_{v \in \bar{\mathcal{N}}(u)} U_v + R + S, \quad (3)$$

- The *unique information* U_v for all $v \in \bar{\mathcal{N}}(u)$ corresponds to the information a neighbor carries independently and no other neighbor has,
- The *redundant information* R is the information that can be found overlapping in two or more neighbors and
- The *synergistic information* S expresses the information that can be captured only if we take the interactions among neighbors into account.

In Figure 1, we provide an illustration of the PID framework. To exemplify this concept we discuss it in the context of the much-used Cora dataset, for which node feature vectors contain binary indication of the presence or absence of certain key words in the abstracts of scientific publications [Sen et al., 2008]. For this dataset unique information takes the form of key words, which are present in only one abstract in a given neighborhood, redundant information refers to keywords, which are repeatedly present without their total number of appearances being of consequence to our learning task, and synergistic information refers to insight that can be gained by observing a certain combination of key words.

3.2 Information Captured by Aggregators

To better understand the information captured by standard GNNs, we first analyze the contribution of each neighboring node to the aggregated representation of a central node.

We assume the structure of an MPNN model, in which each node updates its hidden representation by aggregating information of its neighbors. Further, we denote the message that a given central node u receives from a neighboring node v by $c_{uv} \in \mathbb{R}^d$. Then, c_{uv} can be interpreted as the *contribution* of node v to the hidden state of u and the aggregated messages in (1) can be expressed as $m_u = \sum_{v \in \mathcal{N}(u)} c_{uv}$.

For the Graph Isomorphism Network (GIN) [Xu et al., 2019] and Graph Convolutional Network (GCN) [Kipf and Welling, 2017] we observe that $c_{uv} = f(\mathbf{S}_{uv}, h_v) = \mathbf{S}_{uv} h_v$, where $\mathbf{S} \in \mathbb{R}^{N \times N}$ is a graph shift operator, such as \mathbf{A} or $(\mathbf{D} + \mathbf{I})^{-1/2}(\mathbf{A} + \mathbf{I})(\mathbf{D} + \mathbf{I})^{-1/2}$, and \mathbf{D} denotes the graph’s degree matrix. The contribution of each neighbor v is only determined by its hidden state h_v and the value \mathbf{S}_{uv} of the graph shift operator. For the GAT we observe that $c_{uv} = f(h_u, h_v) = a_{uv} h_v$, where a_{uv} is the attention score that is computed from h_u and h_v . The contribution of each neighbor is also affected by the hidden state of the central node, but is not affected by the other neighbors.

We argue that processing each neighbor individually limits current aggregators, as any interactions among neighbors are ignored by design. Therefore, they can not capture synergistic information between nodes, i.e., the amount of information that is captured equals $\sum_{v \in \mathcal{N}(u)} I(h_u; h_v)$. Consider the example of a neighborhood with two neighbors v_1, v_2 . The information captured by a standard GNN is expressed in terms of the PID as follows, $I(h_u; h_{v_1}) + I(h_u; h_{v_2}) = U_{v_1} + U_{v_2} + 2R$, which is different from the joint mutual information $I(h_u; h_{v_1}, h_{v_2}) = U_{v_1} + U_{v_2} + R + S$. Thus, the captured information from a standard GNN is less than the information present in the neighborhood due to the absence of synergistic information.

To address this problem, we introduce a dependence of the contribution c_{uv} of the neighbor node v on all neighbors of u . Therefore, c_{uv} is now a function not only of h_u and h_v , but also of h_j for $j \in \mathcal{N}(u)$, i.e., $c_{uv} = f(\mathbf{S}, \mathbf{H}_{\mathcal{N}(u)})$. To achieve this, we learn a meaningful ordering of the neighbor nodes using an attention mechanism, and then use an RNN to aggregate the representations of the neighbors.

4 Graph Ordering Attention Layer

We now present the architecture of our *Graph Ordering Attention (GOAT)* layer and highlight its theoretical advantages over other message-passing models. A deep GNN can be constructed by stacking several GOAT layers or combining GOAT layers with other GNN layers. A GOAT layer (illustrated in Figure 2) consists of two parts:

1) The **Ordering Part** (red box in Figure 2) transforms the unordered multiset of neighbor hidden state vectors, each of dimension d , $\{h_1, \dots, h_P\}$, with $P = |\mathcal{N}(u)|$, into an ordered sequence, using an attention mechanism,

$$[h_{\pi(1)}, \dots, h_{\pi(P)}] = \text{OrderingPart}(\{h_1, \dots, h_P\}),$$

where the ordering is given by the permutation function $\pi(\cdot)$.

Specifically, similar to the GAT [Veličković et al., 2018a] model, for each node $v_i \in V$, we first apply a shared linear transformation parameterized by a *weight matrix* $\mathbf{W}_1 \in \mathbb{R}^{d \times d}$ and then perform a shared self-attention mechanism parameterized by $\vec{w}_2 \in \mathbb{R}^{2d}$ to compute the *attention scores*

$$a_{ij} = \text{LeakyReLU}(\vec{w}_2^T [\mathbf{W}_1 h_i \| \mathbf{W}_1 h_j]), \quad (4)$$

for all j such that $v_j \in \mathcal{N}(v_i)$. Then, we sort the coefficients in decreasing order of magnitude

$$a_{i\pi(1)}, \dots, a_{i\pi(P)} = \text{sort}(a_{i1}, \dots, a_{iP}), \quad (5)$$

obtaining a specific permutation π of the nodes in the neighborhood. When all attention scores are different from each other, we observe that the sorting function in (5) is deterministic and permutation invariant. In cases where two or more nodes have equal attention scores, we resort to an additional sorting criterion, described in Appendix A, to ensure that our sorting function is deterministic and permutation invariant.

Once we obtain the permutation π , we construct the sorted sequence of neighbourhood hidden states

$$h_{\text{sorted}(i)} = \left[\frac{e^{a_{i\pi(1)}}}{\sum_{j=1}^Q e^{a_{i\pi(j)}}} \mathbf{W}_1 h_{\pi(1)}, \dots, \frac{e^{a_{i\pi(Q)}}}{\sum_{j=1}^Q e^{a_{i\pi(j)}}} \mathbf{W}_1 h_{\pi(P)} \right]. \quad (6)$$

Note that we use the attention scores to both order the hidden states and, after normalisation via the softmax function, as coefficients for the hidden states. Only due to the occurrence of the attention coefficients in (6) are we able to obtain gradients in the backpropagation algorithm for \mathbf{W}_1 and \vec{w}_2 (the sorting function in (5) is not differentiable). Note further that any self-attention mechanism, such as the GATv2 by Brody et al. [2022], can be used instead of GAT, to obtain the attention scores in a GOAT layer.

2) The **Sequence Modeling Part** (yellow box in Figure 2) takes the ordered sequences of nodes produced by the Ordering Part as input and processes them using an RNN, that is shared across all neighborhoods, to generate the new hidden states. In the PID context, the Bidirectional LSTM [Hochreiter and Schmidhuber, 1997] appears to be the best suited RNN available. Its forget gate allows us to discard redundant information; the input gate is sufficiently expressive to isolate unique information; while its memory states allow for the identification of synergistic information.

$$h_i^{\text{new}} = \text{LSTM}(h_{\text{sorted}(i)}) \in \mathbb{R}^{d_o}. \quad (7)$$

Since we utilize a Bidirectional LSTM the contribution of each node, discussed in Section 3.2, depends on all other hidden states in the neighborhood. Specifically, each direction in the Bidirectional LSTM ensures that both the nodes preceding and succeeding a particular node j are taken into account when calculating the contribution c_{ij} of node j .

Note that the choice of the LSTM is made without loss of generality and any RNN could be chosen in the Sequence Modeling Part. Indeed, in Section 5.3 we will observe results for a variety of RNNs chosen as part of our GOAT layer.

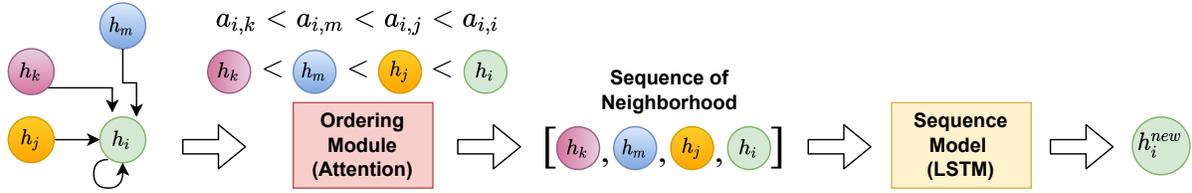


Figure 2: An illustration of the aggregation and update of the representation of node v_i using a GOAT layer. A self-attention mechanism is used in order to obtain a ranking between the nodes of the neighborhood and then the ordered neighborhood is given as input into a sequence model (LSTM) to produce the updated representation of node v_i .

To work with a faster, more scalable implementation we pad all neighborhood sequences with zero vectors to be equal in length to the sequence of hidden states arising in the largest neighborhood in the graph. This allows us to train the LSTM on larger batches of neighborhoods in parallel. The alternative implementation, where neighborhood sequences of different length are fed to the LSTM individually is an equally valid, while slower, implementation.

Multi-Head Attention Ordering. We can also employ multi-head attention to provide additional representational power to our model. We see several advantages in the consideration of multiple heads in our architecture. If only one sensible ordering of the nodes in a neighborhood exists, then multiple heads can help us estimate this ordering more robustly. If on the other hand there exist several sensible orderings of the nodes in a neighborhood, then a multi-head architecture allows us to take all of these into account in our model. Let K be the number of the attention heads. Equation (4) for the k -th attention head is transformed as

$$a_{ij}^k = a^k(\mathbf{W}_1^k h_i, \mathbf{W}_1^k h_j).$$

Then we sort the K sets of attention scores obtaining multiple orderings of the neighborhood, $h_{sorted(i)}^k$ for $k \in \{1, \dots, K\}$. To generate the final representation of the nodes we concatenate the features from the K independent Bidirectional LSTM models, i.e.,

$$h_i^{new} = \parallel_{k=1}^K \text{LSTM}^k \left(h_{sorted(i)}^k \right).$$

Complexity. The time complexity, derived in Appendix B, of a single-head GOAT layer is $\mathcal{O}(|V|d_{\mathcal{O}}d + |E|d_{\mathcal{O}} + |V|d_{\max} \log(d_{\max}) + |V|d_{\max}4d(d_{\mathcal{O}} + d + 3))$, where d_{\max} denotes the maximal degree in the graph. For $d_{\max} \ll d$, the only additional complexity introduced by our model manifests in the multiplicative d_{\max} term in the last summand of the complexity expression. Limiting the maximal degree by applying a sampling strategy, limits the additional complexity our model introduces. Hence, the time complexity of our GOAT model can be comparable to standard MPNNs models. Note that the space complexity of a GOAT layer only exceeds the space complexity of a GAT layer by the space complexity of the LSTM model.

Permutation-Equivariance and Injectivity of GOAT. Recall from Section 2.1 that the permutation-equivariance

and injectivity are desirable properties for a GNN layer to have. We will now prove that our GOAT layer satisfies both of these criteria.

Proposition 4.1 (Permutation-Equivariance of GOAT). Our GOAT layer performs a permutation-equivariant transformation, with respect to node label permutations, of the hidden states corresponding to the nodes in a graph.

The proof of Proposition 4.1 is in Appendix C. Therefore, our GOAT layer is able to benefit from the expressivity of a permutation-sensitive aggregator, while acting on the node’s hidden representations in a permutation-equivariant way.

Our result on the injectivity of a GOAT layer relies on the concept of function approximation in probability. We reproduce the definition of this concept from Hammer [2000].

Definition 4.1. Let \mathcal{X} denote the space of finite lists with elements in \mathbb{R}^q for $q \in \mathbb{N}$ and P be a probability measure on \mathcal{X} . For measurable functions $f_1, f_2 : \mathcal{X} \rightarrow \mathbb{R}^t$ we say that f_1 approximates f_2 with accuracy $\epsilon > 0$ and confidence $\delta > 0$ in probability if $P(x \in \mathcal{X} \mid |f_1(x) - f_2(x)| > \epsilon) < \delta$.

Theorem 1 (Injectivity of GOAT). Assume that for all nodes $u \in V$ the multisets of hidden states corresponding to its neighbors is finite and has elements in \mathbb{R}^q for $q \in \mathbb{N}$. Then, there exists a GOAT layer approximating a measurable function arbitrarily well in probability for which any two distinct multisets are mapped to distinct node representations.

Hence, we have shown that our GOAT layer is sufficiently expressive to approximate any measurable injective function in probability. The proof of Theorem 1 is in Appendix D.

5 Experimental Evaluation

We perform an extensive evaluation of our GOAT model and compare against a wide variety of state-of-the-art GNNs, on three synthetic datasets (in Sections 5.1 and 5.2) as well as on nine node-classification benchmarks (in Section 5.3). Our code is available on [github](#).

Baselines. We compare GOAT against the following state-of-the-art GNNs for node classification: 1) GCN [Kipf and Welling, 2017] the classical graph convolution neural network, 2) GraphSAGE(mean) [Hamilton et al., 2017] that aggregates by taking the elementwise mean value, 3) GraphSAGE(lstm) [Hamilton et al., 2017] that aggregates by feeding the neighborhood hidden states in a random order to an

Table 1: Classification accuracy (\pm standard deviation) on the “Top-2 pooling” synthetic dataset and MSE (\pm standard deviation) results on the synthetic datasets “Betweenness Centrality” and “Effective Size” for two different types of random graphs.

Method	Top-2 pooling (Accuracy)	Betweenness Centrality (MSE)		Effective Size (MSE)	
		N=100, p=0.09	N=1000, p=0.01	N=100, p=0.09	N=1000, p=0.01
GCN	57.35 \pm 4.13	0.0063 \pm 0.0036	0.0020 \pm 0.0008	0.0135 \pm 0.0067	0.00380 \pm 0.00120
GraphSAGE (mean)	61.45 \pm 5.79	0.0401 \pm 0.0158	0.0221 \pm 0.0069	0.0374 \pm 0.0085	0.02430 \pm 0.00560
GraphSAGE (lstm)	65.05 \pm 8.71	0.0094 \pm 0.0073	0.0153 \pm 0.0105	0.0022 \pm 0.0017	0.00080 \pm 0.00020
GIN	56.40 \pm 5.26	0.0083 \pm 0.0052	0.0042 \pm 0.0015	0.0024 \pm 0.0016	0.00070 \pm 0.00030
GAT	53.34 \pm 2.43	0.0409 \pm 0.0158	0.0220 \pm 0.0068	0.0382 \pm 0.0079	0.02480 \pm 0.00560
PNA	61.50 \pm 10.9	0.0115 \pm 0.0089	0.0020 \pm 0.0008	0.0121 \pm 0.0119	0.00137 \pm 0.00035
GOAT(lstm)	69.21 \pm 5.10	0.0038 \pm 0.0019	0.0006 \pm 0.0002	0.0016 \pm 0.0008	0.00020 \pm 0.00008

LSTM, 4) GIN [Xu et al., 2019] the injective summation aggregator, 5) GAT [Veličković et al., 2018a] that aggregates with a learnable weighted summation operation, 6) PNA [Corso et al., 2020] that combines multiple aggregators with degree-scalers. We also compare with 7) a standard MLP that only uses node features and does not incorporate the graph structure. To better understand the impact of the choice of RNN in the GOAT architecture, we provide results from three different GOAT architectures, in which the standard RNN, GRU [Cho et al., 2014] and LSTM are used.

Setup. For a fair comparison we use the same training process for all models adopted by Veličković et al. [2018a]. We use the Adam optimizer [Kingma and Ba, 2015] with an initial learning rate of 0.005 and early stopping for all models and datasets. We perform a hyperparameter search for all models on a validation set. The hyperparameters include the size of hidden dimensions, dropout, and number of attention heads for GAT and GOAT. We fix the number of layers to 2. In our experiments we combine our GOAT layer with a GAT or GCN layer to form a 2-layer architecture. More information about the datasets, training procedure, and hyperparameters of the models are in Appendix E.

5.1 Top-2-Pooling

In this task, we sample Erdős–Rényi random graphs with 1000 nodes and a probability of edge creation of 0.01. We draw 1-dimensional node features from a Gaussian Mixture model with three equally weighted components with means 1, 1 and 2 and standard deviations 1, 4 and 1. Then, we label each node with a function $\phi(\cdot, \cdot)$ of the two 1- or 2-hop neighbors that have the two different largest features, i.e., to each node $u \in V$ we assign a label $y_u = \phi(x_a, x_b)$, where x_a and x_b are the largest, distinct node features of all nodes in the 2-hop neighborhood of u with nodes features at a distance of 2 being down-weighted by a factor of 0.8. We set ϕ to be $\phi(x_a, x_b) = \sqrt{\exp(x_a) + \exp(x_b)}$. Finally, to transform this task to node classification we bin the y values into two equally large classes. We use 60/20/20 percent of nodes for training, validation and testing.

We report the average classification accuracy and the standard deviation across 10 different random graphs in Table 1. *Our GOAT model outperforms the other GNNs with a large margin.* Specifically, our model leads to an 18.36% increase in performance over GAT and 6.65% increase over GraphSage(lstm). In the context of this simulation study, we explain this performance gap with the following hypothesis. To find

the largest element of a set one must consider 2-tuple relationships therefore synergistic information is crucial for this task. An LSTM can easily perform the necessary comparisons with a 2-dimensional hidden space. As nodes are processed they can either be discarded via the forget gate, if they are smaller than the current hidden state, or the hidden state is updated to contain the new feature node. In contrast, typical GNNs need exponentially large hidden dimensions in order to capture the necessary information as they cannot efficiently discard redundant information. We observe that GraphSage(lstm) is the second-best performing model due to its LSTM aggregator. However, it does not learn a meaningful ordering of the nodes that simplifies this task.

5.2 Prediction of Graph Structural Properties

The experiments in this section establish the ability of our GOAT model to predict structural properties of nodes. The first task is to predict the *betweenness centrality* of each node and the second task is to predict the *effective size* of each node. Both of these metrics, defined in Appendix E, are affected by the interactions between the neighbor nodes, so synergistic information is crucial for these tasks. We set the input features to the identity matrix, i.e., $\mathbf{X} = \mathbf{I}$ and use two parameter settings to sample Erdős–Rényi random graphs, namely $(N, p) \in \{(100, 0.09), (1000, 0.1)\}$, where N is the number of nodes and p is the probability of edge creation. We use 60% of nodes for training, 20% for validation and 20% for testing. We train the models by minimizing the Mean Squared Error (MSE).

We report the mean and standard deviation accuracy and MSE across 10 graphs of each type in Table 1. Our model outperforms all models in both tasks and in both graph parameter settings. GOAT can capture the synergistic information between the nodes, which is crucial for predicting the betweenness centrality and effective size. The other aggregators miss the structural information of nodes in neighborhoods. We observe that GraphSAGE(lstm) that uses a random node ordering is not on par with GOAT, indicating that the learned ordering in GOAT is valuable here also.

5.3 Node Classification Benchmarks

We utilize nine well-known node classification benchmarks to validate our proposed model in real-world scenarios originating from a variety of different applications. Specifically, we use 3 citation network benchmark datasets: Cora, CiteSeer [Sen et al., 2008], ogbn-arxiv [Hu et al., 2020],

Table 2: Node classification accuracy using different train/validation/test splits. We highlight the best performing model and underline the second best. Since there exists a single standardised split for Cora and CiteSeer no standard deviations are given.

Method	Cora	CiteSeer	Disease	LastFM Asia	Computers	Photo	CS	Physics
MLP	43.8	52.9	79.10 \pm 0.97	72.27 \pm 1.00	79.53 \pm 0.66	87.89 \pm 1.04	93.76 \pm 0.26	95.85 \pm 0.20
GCN	81.4	67.5	88.98 \pm 2.21	<u>83.58</u> \pm 0.93	90.72 \pm 0.50	93.99 \pm 0.42	92.96 \pm 0.32	96.27 \pm 0.22
GraphSAGE (mean)	77.2	65.3	88.79 \pm 1.95	83.07 \pm 1.19	<u>91.47</u> \pm 0.37	94.32 \pm 0.46	<u>94.11</u> \pm 0.30	96.31 \pm 0.22
GraphSAGE (lstm)	74.1	59.9	90.50 \pm 2.15	86.85 \pm 1.07	91.26 \pm 0.51	94.32 \pm 0.64	93.46 \pm 0.29	96.40 \pm 0.16
GIN	75.5	62.1	90.20 \pm 2.23	82.94 \pm 1.25	84.68 \pm 2.33	90.07 \pm 1.19	92.38 \pm 0.38	96.38 \pm 0.16
GAT	83.0	69.3	89.13 \pm 2.22	77.57 \pm 1.82	85.41 \pm 2.95	90.30 \pm 1.76	92.78 \pm 0.27	96.17 \pm 0.18
PNA	76.4	58.9	86.84 \pm 1.89	83.24 \pm 1.10	90.80 \pm 0.51	<u>94.35</u> \pm 0.68	91.83 \pm 0.33	96.25 \pm 0.21
GOAT(lstm)	84.9	<u>69.5</u>	92.11 \pm 1.88	83.29 \pm 0.91	91.34 \pm 0.50	94.38 \pm 0.66	94.21 \pm 0.42	96.69 \pm 0.31
GOAT(gru)	83.5	70.0	<u>91.97</u> \pm 1.90	83.35 \pm 0.91	91.54 \pm 0.48	94.22 \pm 0.58	93.62 \pm 0.22	96.32 \pm 0.24
GOAT(rnn)	<u>84.2</u>	67.9	91.67 \pm 1.69	83.21 \pm 0.98	89.10 \pm 0.51	92.45 \pm 0.60	93.48 \pm 0.19	<u>96.44</u> \pm 0.20

Table 3: Node classification accuracy on the ogbn-arxiv dataset. We used the same setup and the reported results from Kim and Oh [2022].

Method	ogbn-arxiv
GCN	33.3 \pm 1.2
GraphSAGE	54.6 \pm 0.3
GAT	54.1 \pm 0.5
GOAT(lstm)	55.1 \pm 0.4

1 disease spreading model: Disease [Chami et al., 2019], 1 social network: LastFM Asia [Rozemberczki and Sarkar, 2020], 2 co-purchase graphs: Amazon Computers, Amazon Photo [Shchur et al., 2019] and 2 co-authorship graphs: Coauthor CS, Coauthor Physics [Shchur et al., 2019]. For the GOAT results on the ogbn-arxiv dataset we randomly sample 100 neighbors per node to represent the neighborhoods for faster computation. We report the classification accuracy results in Tables 2 and 3. Our model outperforms the others in eight of nine datasets. This demonstrates the ability of GOAT to capture the interactions of nodes, that are crucial for real-world learning tasks.

5.4 Ablation Studies on the Learned Ordering.

In our GOAT architecture we make the implicit assumption that ordering neighborhoods by the magnitude of the trainable attention scores is an ordering that results in a well-performing model. We now perform several ablation studies where we compare the GOAT model to models with fixed neighborhood node orderings (GOAT-fixed).

Setup. We train our GOAT model on the Cora and Disease datasets, using 8 attention heads and the LSTM aggregator. We store the ordering of the nodes in each neighborhood for each attention head for each epoch. Then, we train various (GOAT-fixed) models that use different fixed orderings extracted from the initial model. Specifically, we train 4 different (GOAT-fixed) models with orderings extracted from the 0, 100, 200, 500 epochs respectively. We run the experiment 3 times and report the results in Table 4.

Discussion. We observe that GOAT-fixed-0, which uses a random ordering, since the ordering is extracted before training, achieves the worst performance. This highlights the importance of a meaningful ordering of the nodes, and the ability of our model to learn one. We also observe that the fixed

Table 4: Accuracy of the GOAT model using fixed orderings extracted from different epochs of the baseline model’s training.

Method	Cora	Disease
GOAT	83.36 \pm 0.42	90.64 \pm 0.40
GOAT-fixed-0	81.56 \pm 0.19	90.48 \pm 0.29
GOAT-fixed-100	83.80 \pm 1.14	90.90 \pm 0.26
GOAT-fixed-200	82.13 \pm 0.27	90.57 \pm 0.59
GOAT-fixed-500	82.27 \pm 0.34	90.50 \pm 0.49

ordering extracted from epoch 100 outperforms the GOAT model. We believe that this phenomenon is associated with the training dynamics of our model. Having a fixed ordering may lead to more stability in the learning of high-performing model parameters not associated with the ordering. For practitioners, learning an ordering in a first run of our model and then training with an extracted fixed ordering may therefore be most advisable.

Additional Ablation Studies In Appendix F.1 we investigate the potential use of the GATv2 model instead of the GAT model in a GOAT layer and find that the two model variants perform comparably. In Appendix F.2 we observe the GOAT model to significantly outperform the Janosy Pooling approach on the Cora, CiteSeer and Disease datasets. In Appendix F.3 we find the optimal number of attention heads in a GOAT layer to be related to complexity of the learning task on a given dataset. In particular, we observe one attention head to yield optimal performance on Cora, four attention heads are optimal for CiteSeer, while eight attention heads resulted in the best performing model in the Disease dataset.

6 Conclusion

We have introduced a novel view of learning on graphs by introducing the Partial Information Decomposition to the graph context. This has allowed us to identify that current aggregation functions used in GNNs often fail to capture synergistic and redundant information present in neighborhoods. To address this issue we propose the Graph Ordering Attention (GOAT) layer, which makes use of a permutation-sensitive aggregator capable of capturing synergistic and redundant information, while maintaining the permutation-equivariance property. The GOAT layer is implemented by first learning an ordering of nodes using a self-attention and by then apply-

ing an RNN to the ordered representations. This theoretically grounded architecture yields improved accuracy in the node classification and regression tasks on both synthetic and real-world networks.

References

- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Trans. Neural Netw.*, 20(1):61–80, 2009. doi: 10.1109/TNN.2008.2005605.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv:2104.13478*, 2021.
- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *ICML*, pages 1263–1272, 2017.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *ICLR*, 2019.
- Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn: Geometric graph convolutional networks. In *ICLR*, 2020.
- Ryan L. Murphy, Balasubramaniam Srinivasan, Vinayak Rao, and Bruno Ribeiro. Janossy pooling: Learning deep permutation-invariant functions for variable-size inputs. In *ICLR*, 2019.
- Edward Wagstaff, Fabian B. Fuchs, Martin Engelcke, Michael A. Osborne, and Ingmar Posner. Universal approximation of functions on sets. *arXiv:2107.01959*, 2021.
- Emilio Bizzi and Vincent C. K. Cheung. The neural origin of muscle synergies. *Front. Comput. Neurosci.*, 7: 51, 2013. ISSN 1662-5188. doi: 10.3389/fncom.2013.00051. URL <https://www.frontiersin.org/article/10.3389/fncom.2013.00051>.
- Naama Brenner, Steven Strong, Roland Koberle, William Bialek, and R Steveninck. Synergy in a neural code. *Neural computation*, 12:1531–52, 08 2000. doi: 10.1162/089976600300015259.
- José Manuel Pérez-Pérez, Héctor Candela, and José Luis Micol. Understanding synergy in genetic interactions. *Trends in Genetics*, 25(8):368–376, 2009. ISSN 0168-9525. doi: <https://doi.org/10.1016/j.tig.2009.06.004>. URL <https://www.sciencedirect.com/science/article/pii/S0168952509001322>.
- Paul L. Williams and Randall D. Beer. Nonnegative decomposition of multivariate information. *arXiv:1004.2515*, 2010. URL <http://arxiv.org/abs/1004.2515>.
- Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. Order matters: Sequence to sequence for sets. In *ICLR*, 2016. URL <http://arxiv.org/abs/1511.06391>.
- William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *NIPS*, 2017.
- M. Gori, G. Monfardini, and F. Scarselli. A new model for learning in graph domains. In *IJCNN*, pages 729–734, 2005. doi: 10.1109/IJCNN.2005.1555942.
- Edward Wagstaff, Fabian B. Fuchs, Martin Engelcke, Ingmar Posner, and Michael Osborne. On the limitations of representing functions on sets. In *ICML*, 2019.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *ICLR*, 2018a.
- Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks? In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=F72ximsx7C1>.
- Shuo Zhang and Lei Xie. Improving attention mechanism in graph neural networks via cardinality preservation. In *IJCAI*, 2020. ISBN 9780999241165. doi: 10.24963/ijcai.2020/194. URL <http://dx.doi.org/10.24963/ijcai.2020/194>.
- Mathias Niepert, Mohamed Ahmed, and Konstantin Kutikov. Learning convolutional neural networks for graphs. In *International conference on machine learning (ICML)*, pages 2014–2023. PMLR, 2016.
- Hongyang Gao, Zhengyang Wang, and Shuiwang Ji. Large-scale learnable graph convolutional networks. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1416–1424, 2018.
- Petar Veličković, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. Deep graph infomax. In *ICLR*, 2018b.
- Zhen Peng, Wenbing Huang, Minnan Luo, Qinghua Zheng, Yu Rong, Tingyang Xu, and Junzhou Huang. Graph representation learning via graphical mutual information maximization. In *WWW*, 2020. URL <https://arxiv.org/abs/2002.01169>.
- Gongxu Luo, Jianxin Li, Hao Peng, Carl Yang, Lichao Sun, Philip S. Yu, and Lifang He. Graph entropy guided node embedding dimension selection for graph neural networks. *arXiv:2105.03178*, 2021. URL <https://arxiv.org/abs/2105.03178>.
- George Dasoulas, Giannis Nikolentzos, Kevin Scaman, Aladin Virmaux, and Michalis Vazirgiannis. Ego-based entropy measures for structural representations. *ICASSP*, 2020. URL <https://arxiv.org/abs/2003.00553>.
- Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI Mag.*, 29(3):93, 2008. doi: 10.1609/aimag.v29i3.2157. URL <https://ojs.aaai.org/index.php/aimagazine/article/view/2157>.
- Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Comput.*, 9(8):1735–1780, 11 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <https://doi.org/10.1162/neco.1997.9.8.1735>.
- Barbara Hammer. On the approximation capability of recurrent neural networks. *Neurocomputing*, 31(1-4):107–123, 2000.

- Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. In *Advances in Neural Information Processing Systems*, 2020.
- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, pages 103–111, 2014.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *arXiv:2005.00687*, 2020.
- Ines Chami, Rex Ying, Christopher Ré, and Jure Leskovec. Hyperbolic graph convolutional neural networks. *NeurIPS*, 2019.
- Benedek Rozemberczki and Rik Sarkar. Characteristic Functions on Graphs: Birds of a Feather, from Statistical Descriptors to Parametric Models. In *CIKM*, page 1325–1334. ACM, 2020.
- Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. In *R2L Workshop at NeurIPS*, 2019.
- Dongkwan Kim and Alice Oh. How to find your friendly neighborhood: Graph attention design with self-supervision, 2022. URL <https://arxiv.org/abs/2204.04879>.
- Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric, 2019. URL <https://arxiv.org/abs/1903.02428>.
- Roy M Anderson and Robert M May. *Infectious diseases of humans: dynamics and control*. Oxford university press, 1992.
- K-I Goh, Eulsik Oh, Byungnam Kahng, and Doochul Kim. Betweenness centrality correlation in social networks. *Phys. Rev. E*, 67(1):017101, 2003.
- Martin Everett and Stephen Borgatti. Unpacking burt’s constraint measure. *Soc. Netw.*, 62:50–57, 07 2020. doi: 10.1016/j.socnet.2020.02.001.

Appendix of Graph Ordering Attention Networks

A Resolving the Ties in Attention Scores

When several neighboring nodes have equal attention scores in a given neighborhood, then a simple ordering by attention scores is not deterministic and permutation-invariant. Therefore, we introduce an additional sorting criterion to resolve the ties between equal attention score nodes. Specifically, in this additional sorting criterion we compare the hidden state elements of nodes successively until we detect unequal elements, which are then used to order the nodes in ascending order by discriminating element. In Algorithm 1 we lay out the necessary steps in pseudocode.

For example, assume nodes v_j, v_k in the neighborhood of central node v_i are assigned equal attention scores, but have different embedding vectors $h_j, h_k \in \mathbb{R}^d, h_j \neq h_k$,

$$\begin{aligned} h_j &= [h_{j1}, h_{j2}, \dots, h_{jd}], \\ h_k &= [h_{k1}, h_{k2}, \dots, h_{kd}]. \end{aligned}$$

Then, there exists at least one index $\ell \in \{1, \dots, d\}$ such that $h_{j\ell} \neq h_{k\ell}$, since $h_j \neq h_k$. We consider the smallest index ℓ for which $h_{j\ell} \neq h_{k\ell}$. If $h_{j\ell} < h_{k\ell}$ then we put v_j first in our learned ordering, otherwise we put v_k first. To illustrate this further, if $h_j = [1, 2, 3]$ and $h_k = [1, 3, 3]$ and v_j, v_k have equal attention scores, then v_j will be put ahead of v_k in our ordering, since $h_{j2} = 2 < 3 = h_{k2}$. This sorting criterion also resolves ties of more than two nodes. Therefore, our sorting function is deterministic and permutation-invariant, even for the case of equal attention scores between two or more nodes.

Algorithm 1 Resolve Ties

```
1: Input: Nodes  $v_j, v_k$  with embedding vectors  $h_j, h_k \in \mathbb{R}^d$  and equal attention scores  $a_{ij}, a_{ik}, a_{ij} = a_{ik}$ .
2: for  $q = 1$  to  $d$  do
3:   if  $h_{jq} < h_{kq}$  then
4:     first  $\leftarrow v_j$ 
5:     second  $\leftarrow v_k$ 
6:     break
7:   else if  $h_{jq} > h_{kq}$  then
8:     first  $\leftarrow v_k$ 
9:     second  $\leftarrow v_j$ 
10:    break
11:   end if
12: end for
13: Output: first, second
```

B Computational Complexity

Our GOAT model requires the computation of the three following steps:

1. **Computation of attention scores.** The computational complexity of the GAT or GATv2 model, used to calculate the attention scores, is $\mathcal{O}(|V|d_O d + |E|d_O)$, where d is input dimensions and d_O is the output dimensions [Brody et al., 2022].

2. **Sorting attention scores.** The sorting operation requires $\mathcal{O}(d_u \log(d_u))$ steps for each node u , the degree of which we denote by d_u . To parallelize the computation across the nodes, we pad all the neighborhoods to the maximum degree in the graph d_{\max} . Therefore, the complexity of the second step is $\mathcal{O}(|V|d_{\max} \log(d_{\max}))$.

3. **Update hidden states using an RNN.** The total number of parameters in a standard LSTM network is equal to $W = 4h(d_O + h) + 4h$, where d_O is the number of input units, and h is the number of output units. The evaluation of the 5 activation functions and 3 element-wise products involves a total complexity of $\mathcal{O}(8h)$. Therefore, the computational complexity per time step is $\mathcal{O}(4h(d_O + h + 3))$. So, the complexity of computing the representation of each node is equal to $\mathcal{O}(d_u 4h(d_O + h + 3))$. Since we parallelize the computation across the graphs by using the padded sequence, we end up with complexity $\mathcal{O}(|V|d_{\max} 4h(d_O + h + 3))$.

Therefore, the final complexity of our model is

$$\begin{aligned} \mathcal{O}(|V|d_O d + |E|d_O + |V|d_{\max} \log(d_{\max}) \\ + |V|d_{\max} 4h(d_O + h + 3)). \end{aligned}$$

If the maximal degree d_{\max} of the graph is large, we can apply a neighbor sampling strategy like in the GraphSAGE model, instead of working with the whole neighborhood. This allows us to assume that $d_{\max} \ll d, d2$. In this case our complexity is $\mathcal{O}(|V|d_O d + |E|d_O + |V|W d_{\max})$. Limiting the maximal degree, limits the additional complexity our model introduces.

C Proof of Proposition 4.1

Typically GNNs construct permutation-equivariant functions on graphs by applying a permutation-invariant local function over the neighborhood of each node [Bronstein et al., 2021]. To establish the permutation-equivariance of the GOAT layer it therefore suffices to show that the node-wise operation performed by our GOAT layer is permutation-invariant. To do so we make use of the following proposition which concerns the permutation-invariance of composed functions.

Proposition C.1. For any function $f : X \rightarrow Y$ and for any permutation-invariant function $g : Z \rightarrow X$, their composition $f \circ g$ is permutation-invariant.

Now since the GOAT layer is formed by the composition of the Sequence Modelling Part and the Ordering Part, by Proposition C.1 it suffices to show that the Ordering Part is permutation-invariant to establish the permutation-invariance of their composition. Recall, that in the Ordering Part of the GOAT layer we implement an attention mechanism on the hidden states of the central node and each neighboring node. Then, nodes are ordered according to the magnitude of the attention scores. Crucially, these computations are independent of the node labelling, even when equal attention scores arise as we show in Appendix A, making the Ordering Part of the GOAT layer permutation-invariant. Consequently, we apply a local permutation-invariant function rendering the action of the GOAT layer on the graph permutation-equivariant.

Table 5: Summary of the datasets used in our experiments.

	Cora	CiteSeer	Disease	LastFM Asia	Computers	Photo	CS	Physics
# Nodes	2708	3327	1044	7624	13752	7650	18333	34493
# Edges	5429	4732	1043	27806	491722	238162	163788	495924
# Features/Node	1433	3703	1000	128	767	745	6805	8415
# Classes	7	6	2	18	10	8	10	8
# Training Nodes	140	120	312	4574	9625	5354	12832	24144
# Validation Nodes	300	500	105	1525	1376	765	1834	3450
# Test Nodes	1000	1000	627	1525	2751	1531	3667	6899

D Proof of Theorem 1

Our GOAT layer is a functional composition of the Ordering Part and the Sequence Modeling Part described in Section 4. Since the composition of two injective functions is injective itself, it suffices to show that each of the two components is injective.

We begin by considering the Ordering Part, which maps a multiset of hidden states to an ordered multiset of hidden states leaving the elements of these multisets unchanged. If therefore, for two multisets the same output is generated in the Ordering Part, then their elements are equal. Two multisets with all equal elements are equal themselves. Therefore, the Ordering Part of our GOAT layer is an injective function.

We now consider the Sequence Modeling Part. Following [Hammer, 2000, p. 3] we refer to an activation function σ as a *squashing activation function* if $\sigma : \mathbb{R} \rightarrow [0, 1]$ is monotonous with $\lim_{x \rightarrow -\infty} \sigma(x) = 0$ and $\lim_{x \rightarrow \infty} \sigma(x) = 1$. In the now following Theorem 2 we combine and rephrase the universal approximation results in Theorem 3 and Corollary 4 in [Hammer, 2000, pp. 6, 8].

Theorem 2. Any measurable function $f : \mathcal{X} \rightarrow \mathbb{R}^t$ can be approximated arbitrarily well in probability by a recurrent neural network $\phi \circ \tilde{\rho}_y$, where $\rho : \mathbb{R}^{q+1} \rightarrow \mathbb{R}$ is a feedforward neural network without a hidden layer and either a squashing activation function or a locally Riemann integrable and non polynomial activation function. $h : \mathbb{R} \rightarrow \mathbb{R}^t$ is either a linear mapping or a multilayer network with one hidden layer with squashing or locally Riemann integrable and nonpolynomial activation function and linear outputs.

Theorem 2 applies to recurrent neural networks which are formally defined by [Hammer, 2000, p. 3] to be composed of two feedforward neural networks $\rho : \mathbb{R}^{q+1} \rightarrow \mathbb{R}$ and $\phi : \mathbb{R} \rightarrow \mathbb{R}^t$. The function ρ is referred to as the “recursive part” since it gives rise to the function

$$\tilde{\rho}_y([h_1, \dots, h_i]) = \begin{cases} y & \text{for } i = 0; \\ \rho(h_i, \tilde{\rho}_y([h_1, \dots, h_{i-1}])) & \text{otherwise.} \end{cases}$$

The function ϕ is simply referred to as the “feedforward part” since it succeeds the recursive part and plays the role of a standard readout function in a deep learning architecture. That is to say, that a recurrent neural network f is defined to equal the composition of $f = \phi \circ \tilde{\rho}_y$. Now all three of the RNNs we consider as possible aggregators, the standard RNN, LSTM and GRU, satisfy the conditions of this formal

definition of recurrent neural networks and therefore, Theorem 2 can be applied to our three considered RNNs.

Now that we have shown that Theorem 2 applies to all RNNs we consider in the Sequence Modeling Part, we can use it to establish that the function learned in our Sequence Modeling Part can approximate any measurable function arbitrarily well in probability. This notably includes all measurable injective functions, thus providing us with the desired result.

E Experimental Details

Included in this supplementary material is our implementation, which is built upon the open source library *PyTorch Geometric* (PyG) under MIT license [Fey and Lenssen, 2019]. The experiments are run on an Intel(R) Xeon(R) CPU E5-1607 v2 @ 3.00GHz processor with 128GB RAM and a NVIDIA Corporation GP102 TITAN X GPU with 12GB RAM.

Datasets Details. In our experiments we utilize nine well-known node classification benchmarks. We describe them below:

- 2 citation network benchmark datasets: Cora, CiteSeer [Sen et al., 2008], where nodes represent scientific papers, edges are citations between them, and node labels are academic topics. We follow the experimental setup of Kipf and Welling [2017] and use 140 nodes for training, 300 for validation and 1000 for testing. We optimize hyperparameters on Cora and use the same hyperparameters for CiteSeer.
- 1 disease spreading model: Disease [Chami et al., 2019]. It simulates the SIR disease spreading model [Anderson and May, 1992], where the label of a node indicates if it is infected or not. We follow the experimental setup of Chami et al. [2019] and use 30/10/60% for training, validation and test sets and report the average results from 10 different splits.
- 1 social network: LastFM Asia [Rozemberczki and Sarkar, 2020]. Nodes are LastFM users from Asian countries and edges are mutual follower relationships between them. The label of each node is the country of the user. We use 60/20/20% for training, validation and test sets and report the average results from 10 different splits.

Table 6: Comparison between different attention mechanisms in the GOAT layer. We report the classification accuracy (\pm standard deviation) on the Cora, CiteSeer, Disease and “Top-2 pooling” datasets and MSE (\pm standard deviation) results on the synthetic datasets “Betweenness Centrality” and “Effective Size” for two different types of random graphs.

Method	Cora	CiteSeer	Disease	Top-2 pooling	Betweenness Centrality		Effective Size	
					(100,0.09)	(1000,0.01)	(100,0.09)	(1000,0.01)
GOAT(gat)	84.9	69.5	92.11 ± 1.88	69.21 ± 5.10	0.0038 ± 0.0019	0.0006 ± 0.0002	0.0016 ± 0.0008	0.0002 ± 0.000082
GOAT(gatv2)	83.1	69.3	91.28 ± 1.75	67.34 ± 5.24	0.0038 ± 0.0022	0.0006 ± 0.0001	0.0013 ± 0.0008	0.0001 ± 0.000037

- 2 co-purchase graphs: Amazon Computers, Amazon Photo [Shchur et al., 2019]. Nodes represent products and edges represent that two products are frequently bought together. The node label indicates the product category. We use 70/10/20% for training, validation and test sets and report the average results from 10 different splits. We optimize hyperparameters on Computers and use the same hyperparameters for Photo.
- 2 co-authorship graphs: Coauthor CS, Coauthor Physics [Shchur et al., 2019]. Nodes represent authors that are connected by an edge if they co-authored a paper. Given paper keywords for each author’s papers as node features, the task is to identify the field of study of the authors authors. We use 70/10/20% for training, validation and test sets and report the average results from 10 different splits.
- 1 large citation network: ogbn-arxiv [Hu et al., 2020]. Each node is an arXiv paper and an edge indicates that one paper cites another one. The task is to predict the 40 subject areas of arXiv CS papers. We use the public split by publication dates provided by the original paper.

We report further summary statistics of these datasets in Table 5.

In our synthetic tasks we predic the betweenness centrality $b(u)$ and the effective size $e(u)$.

The betweenness centrality $b(u)$ is a measure of centrality of a node u based on shortest paths involving u . It has many applications in network science, as it is a useful metric for analyzing communication dynamics [Goh et al., 2003]. It can be computed using the following equation

$$b(u) = \sum_{s,t \in V} \frac{\sigma(s,t|u)}{\sigma(s,t)},$$

where $\sigma(s,t)$ is the number of distinct shortest paths between vertices s and t , and $\sigma(s,t|u)$ is the number of these shortest paths passing through u .

The effective size $e(u)$ [Everett and Borgatti, 2020] of node u is based on the concept of redundancy and for the case of unweighted and undirected graphs, can be computed as

$$e(u) = n - \frac{2q}{n},$$

where q is the number of ties in the subgraph induced by the node set $\mathcal{N}(u)$ (excluding ties involving u) and $n = |\mathcal{N}(u)|$ is the number of neighbors (excluding the central node).

Synthetic Experiments: Prediction of Graph Structural Properties (node regression) For the GCN and GraphSage

model we transform the input features with a linear layer and then use 2 convolutional layers followed by 1 linear layer. To optimize the hyper-parameters we perform a grid-search on the following values: linear = {4, 8, 16, 32, 64} for the first linear layer, conv1 = {4, 8, 16, 32, 64} for the first convolutional layer, conv2 = {4, 8, 16, 32} for the second convolutional layer. For the GAT and GOAT model we optimize the following hyper-parameters: nheads = {1, 4, 8} for the number of heads, conv1 = {4, 8, 16, 32, 64} for the first convolutional layer, conv2 = {4, 8, 16, 32, 64} for the second convolutional layer. Also for the GOAT model, we use one GOAT layer and one GCN or GAT layer as the second layer. Specifically, for the “Betweenness Centrality” and “Effective Size” tasks we used GAT as the second layer, and for the “Top-2 pooling” task we used GCN. For the PNA model we optimize the following hyper-parameters: linear = {4, 8, 16, 32, 64} for the first linear layer, conv1 = {4, 8, 16, 32, 64} for the first convolutional layer, conv2 = {4, 8, 16, 32} for the second convolutional layer, aggregators = {‘mean’, ‘min’, ‘max’, ‘std’} for the aggregators scalers = {‘identity’, ‘amplification’, ‘attenuation’, ‘linear’} for the scalers. We search for the best model on ($N = 100, p = 0.09$) and we use the same models for the other configuration of each task ($N = 1000, p = 0.1$).

Node classification Benchmarks. For node classification benchmarks we follow the same model configurations as with node regression above and we just remove the last linear layers from all the models. For the GOAT model, we used GAT as a second layer in Cora, CiteSeer and Disease and GCN as a second layer in LastFM Asia, Computers, Photo, CS and Physics datasets.

F Additional Experiments

F.1 Comparison of GAT and GATv2

We have performed additional experiments investigating the GATv2 attention mechanism as part of a GOAT layer and found it to yield comparable performance to the original GAT attention mechanism. We follow the same setup as the main experiments and we use an LSTM to aggregate the hidden states of the ordered neighbors. The results are reported in Table 6.

F.2 Comparison with Janossy Pooling

We furthermore investigated how our GOAT model compares to the Janossy Pooling model. In Table 7 we observe the GOAT model to significantly outperform three different hyperparametrizations of the Janossy Pooling model on the Cora, CiteSeer and Disease datasets.

Table 7: Comparison of the accuracies attained by our GOAT architecture and three different hyperparametrizations, given in the format (k_1, k_2) , of the Janossy Pooling model[Murphy et al., 2019].

Method	Cora	CiteSeer	Disease
Janossy Pooling(5,5)	79.0	64.2	87.21 \pm 1.93
Janossy Pooling(15,5)	80.8	65.8	87.15 \pm 1.86
Janossy Pooling(20,20)	80.2	64.7	87.19 \pm 1.94
GOAT(lstm)	84.9	<u>69.5</u>	92.11 \pm 1.88
GOAT(gru)	83.5	70.0	91.97 \pm 1.90
GOAT(rnn)	<u>84.2</u>	67.9	91.67 \pm 1.69

Table 8: Accuracy scores of the GOAT architecture when the number of attention heads is varied.

Method	Cora	CiteSeer	Disease
GOAT(lstm)_{1h}	84.9	67.9	89.14 \pm 2.99
GOAT(lstm)_{2h}	83.1	68.2	90.78 \pm 1.93
GOAT(lstm)_{4h}	82.8	69.5	91.32 \pm 2.71
GOAT(lstm)_{8h}	82.9	68.8	92.11 \pm 1.88

F.3 Varying the Number of Attention Heads in GOAT

In this experiment, we examine the performance of our model using different numbers of attention heads. We use the standard configuration of our GOAT model, i.e., a GAT attention mechanism paired with an LSTM aggregator. We report the results in Table 8.

We observe that for datasets, for which the learning tasks are known to be relatively simple, i.e., the Cora dataset, a small number of attention heads is sufficient to achieve the best performance. For datasets with complex node interactions and high amount of synergistic information, learning a large number of neighborhood orderings, appears to be beneficial.