# Exact Computation of Any-Order Shapley Interactions for Graph Neural Networks

## Maximilian Muschalik\*, Fabian Fumagalli\*, Paolo Frazzetto, Janine Strotherm, Luca Hermes, Alessandro Sperduti, Eyke Hüllermeier, and Barbara Hammer



## Background

A Graph Neural Network (GNN) makes prediction for a graph g:  $\sigma$ : output layer  $f_g(\mathbf{X}) := \sigma(\Psi(\{\{f_i(\mathbf{X}) \mid v_i \in V\}\})) \quad \text{with}$ Möbius Interactions  $m: \mathcal{P}(N) \to \mathbb{R}$  are the basis of explanations:

 $m(S) := \sum (-1)^{|S| - |T|} \nu(T) \quad \text{and they recover} \quad \nu(T) = \sum m(S)$ 

Shapley Values and Interactions summarize/aggregate the Möbius interactions into lower-order explanations.

> Shapley Values measure feature attribution of nodes > Shapley Interactions measure synergy between nodes

## References

### For full references, see our accompanying paper.

Shapley, L. S. (1953). A Value for n-Person Games. Contributions to the Theory of Games, Volume II, pages 307–318. Grabisch and Roubens (1999). An Axiomatic Approach to the Concept of Interaction Among Players in Cooperative Games. Int. J. Game Theory 28(4):547–565

Bordt and von Luxburg (2023). From Shapley Values to Generalized Additive Models and Back. AISTATS'23 Lundberg et al. (2017). A Unified Approach to Interpreting Model Predictions. NeurIPS'17

Fumagalli et al. (2023). SHAP-IQ: Unified Approximation of any-order Shapley Interactions. NeurIPS'23

Ye et al. (2023). SAME: uncovering GNN black box with structure-aware shapley-based multipiece explanations. NeurIPS'23 Amara et al. (2022) Graphframex: Towards systematic evaluation of explainability methods for graph neural networks. LoG'22

## **Contribution Summary**



## **Empirical Results**

