# QGShap: Quantum Acceleration for Faithful GNN Explanations

**Haribandhu Jena**                                        *haribandhu.jena@niser.ac.in*
*School of Computer Sciences*
*National Institute of Science Education and Research*
*An OCC of Homi Bhabha National Institute, India*

**Jyotirmaya Shivottam**                                   *jyotirmaya.shivottam@niser.ac.in*
*School of Computer Sciences*
*National Institute of Science Education and Research*
*An OCC of Homi Bhabha National Institute, India*

**Subhankar Mishra**                                       *smishra@niser.ac.in*
*School of Computer Sciences*
*National Institute of Science Education and Research*
*An OCC of Homi Bhabha National Institute, India*

## Abstract

Graph Neural Networks (GNNs) have become indispensable in critical domains such as drug discovery, social network analysis, and recommendation systems, yet their black-box nature hinders deployment in scenarios requiring transparency and accountability. While Shapley value-based methods offer mathematically principled explanations by quantifying each component's contribution to predictions, computing exact values requires evaluating $2^n$ coalitions (or aggregating over $n!$ permutations), which is intractable for real-world graphs. Existing approximation strategies sacrifice either fidelity or efficiency, limiting their practical utility. We introduce QGSHAP[1], a quantum computing approach that leverages amplitude amplification to achieve quadratic speedups in coalition evaluation while maintaining exact Shapley computation. Unlike classical sampling or surrogate methods, our approach provides fully faithful explanations without approximation trade-offs for tractable graph sizes. We conduct empirical evaluations on synthetic graph datasets, demonstrating that QGSHAP achieves consistently high fidelity and explanation accuracy, matching or exceeding the performance of classical methods across all evaluation metrics. These results collectively demonstrate that QGSHAP not only preserves exact Shapley faithfulness but also delivers interpretable, stable, and structurally consistent explanations that align with the underlying graph reasoning of GNNs. The implementation of QGSHAP is available at https://github.com/smlab-niser/qgshap.

## 1 Introduction

Graph neural networks (GNNs) have gained widespread use for learning from graph-structured data in critical applications such as molecular chemistry (Reiser et al., 2022), social network analysis (Li et al., 2023), and recommendation systems (Wu et al., 2022). They excel at capturing complex relationships and patterns within interconnected data (Joshi & Mishra, 2022), enabling breakthroughs in areas like drug

---

[1]To appear in the CCIS series (Springer Nature), QC+AI Workshop at AAAI 2026.

1

discovery, social network analysis, and recommendation systems. However, despite their success, GNNs often function as 'black boxes', making it difficult for users and stakeholders to understand how they arrive at their decisions (Agarwal et al., 2023). This opacity poses significant challenges in domains, where transparency, trust, and accountability are essential. Additionally, the complexity of GNN architectures (Gilmer et al., 2017; Kipf & Welling, 2017; Veličković et al., 2018) and the diversity of graph data further complicate efforts to interpret their predictions.

Building on recent advances in GNN explainability, researchers have moved beyond node and edge-level explanation methods, such as GNNExplainer (Ying et al., 2019), PGExplainer (Luo et al., 2020), and GraphLIME (Huang et al., 2022), toward approaches that capture more complex structural patterns in graphs. While GNNExplainer and PGExplainer use gradient and perturbation-based techniques to identify important components, they often produce explanations that are faithful but unstable, especially on complex benchmarks (Agarwal et al., 2023). Surrogate learning methods like GraphLIME aim to explain local feature relationships, but they still explain only the node features. More recent work, such as SubgraphX (Yuan et al., 2021), leverages Shapley values (Kuhn & Tucker, 1953) combined with Monte Carlo Tree Search (Silver et al., 2017) to identify entire explanatory subgraphs. This shift enables explanations that are both more faithful and interpretable at a higher semantic level, moving the field toward more robust and meaningful forms of interpretability.

Motivated by this shift toward subgraph-level reasoning, Shapley value-based (Lundberg & Lee, 2017) explainability methods have emerged as the principled foundation underlying such scoring, offering a mathematically rigorous way to quantify how each node or subgraph contributes to a model's prediction (Yuan et al., 2021; Akkas & Azad, 2024; Duval & Malliaros, 2021). By aggregating the marginal impact of each component across all possible coalitions, Shapley values provide fairness and completeness in attribution (Lundberg & Lee, 2017). Yet, the very strength of this formulation - its exhaustive consideration of all $2^n$ combinations, renders it computationally prohibitive for real-world graphs. Classical approaches have sought to approximate these values through sampling, Monte Carlo estimation, or surrogate modeling, but such strategies unavoidably trade off either fidelity or efficiency (Yuan et al., 2021; Duval & Malliaros, 2021; Akkas & Azad, 2024; Muschalik et al., 2025).

To move beyond this bottleneck while retaining Shapley's axiomatic benefits, recent advances in quantum computing have introduced algorithms that exploit amplitude amplification to achieve quadratic speedups for combinatorial evaluation tasks, including subset and coalition scoring (Montanaro, 2015; Burge et al., 2025). Building on these developments, we propose QGSHAP for GNN explainability that leverages amplitude amplification to accelerate coalition evaluation, while maintaining exact Shapley computation. Empirical evaluations on synthetic and small real-world datasets demonstrate that our method achieves exact Shapley faithfulness, suggesting that quantum algorithms can play a transformative role in scaling the explainability of GNNs.

**Contributions.**

1. We introduce QGSHAP, demonstrating that quantum amplitude-estimation techniques can accelerate brute-force Shapley value computation for GNN explanations, achieving faithful, theoretically grounded node attributions with quadratic query speedup over classical approaches.

2. QGSHAP provides exact Shapley-based explanations and surpasses existing explainers on explanation-quality metrics across synthetic graph benchmarks.

## 2 Related Work

GNNs have established themselves as highly effective frameworks for learning and reasoning over complex structured data (Wu et al., 2021). Despite their remarkable predictive capabilities and widespread success across numerous domains, the internal decision-making mechanisms of GNNs often remain largely opaque (Agarwal et al., 2023). Early research in GNN explainability focused on attributing importance to individual nodes, edges, or features through gradient-based techniques such as Saliency Analysis (SA) (Zeiler & Fergus, 2014), Class Activation Mapping (CAM) (Zhou et al., 2016), and Guided Backpropagation

(Guided BP) (Springenberg et al., 2015); decomposition-based methods like Layer-wise Relevance Propagation (LRP) (Schwarzenberg et al., 2019) and Excitation Backpropagation (Excitation BP) (Pope et al., 2019); and perturbation-based approaches that assess model sensitivity to input modifications. Among these, GNNExplainer (Ying et al., 2019) learns soft masks over graph components to maximize the mutual information between original and perturbed predictions, while PGExplainer (Luo et al., 2020) extends this by learning a parameterized probabilistic model that generalizes edge importance prediction across graphs. Surrogate-based frameworks, such as GraphLIME (Huang et al., 2022), adopt locally interpretable linear models to approximate the neighborhood-level decision boundary of GNNs. These explanation methods are locally accurate but unstable and lack broader, human-friendly explanations (Agarwal et al., 2023).

Recent advances in GNN explainability have centered on Shapley value-based methods (Kuhn & Tucker, 1953), each introducing distinct approximations to make computation tractable for real-world graph data. SubgraphX (Yuan et al., 2021) employs Monte Carlo Tree Search (Silver et al., 2017) and approximates Shapley values through limited sampling, which, while efficient for small graphs, becomes impractically slow for larger or denser graphs due to the exponential coalition space. GraphSVX (Duval & Malliaros, 2021) constructs a surrogate model on a perturbed dataset and samples coalitions, but its model-agnostic approach undersamples mid-sized coalitions, potentially reducing explanation fidelity. GNNShap (Akkas & Azad, 2024) leverages GPU parallelism and batching to accelerate Shapley value estimation, achieving significant speedups over prior methods, yet fundamentally remains an approximation technique reliant on sampling rather than exact computation. In contrast, GraphSHAP-IQ (Muschalik et al., 2025) exploits the structure of message-passing GNNs with linear global pooling and output layers to compute exact any-order Shapley interactions, scaling near-linearly for sparse and shallow graphs and requiring far fewer model calls than model-agnostic baselines. However, for deep architectures, densely connected graphs, or when the largest $\ell$-hop neighborhood exceeds a practical threshold, GraphSHAP-IQ must introduce a hyperparameter to limit the highest order of computed interactions, trading exactness for tractability. Critically, its theoretical guarantees break down for GNNs with nonlinear readout functions, as interactions extend beyond receptive fields and the sparse interaction property is lost, substantially increasing complexity. As a result, GraphSHAP-IQ cannot be directly applied to quantum GNNs such as Equivariant Quantum Graph Circuits (Mernyei et al., 2022), where nonlinear quantum interactions and entanglement fundamentally violate the linearity and decomposability assumptions required for the explainability frameworks.

The computation of Shapley values for GNN explanations is intractable in general. Exact computation is #P-complete for classical (explicit) cooperative games (Deng & Papadimitriou, 1994), and becomes $\mathsf{FP}^{\#\mathsf{P}}$-hard (and in some cases #P-complete) in succinct graph or query settings such as RPQs and CRPQs (Khalil & Kimelfeld, 2023; Bienvenu et al., 2024). Thus, the exponential coalition enumeration $2^n$ should be viewed as a symptom rather than the formal cause. Consequently, precise evaluation involves considering all $2^n$ possible node subsets for a graph of $n$ nodes in the worst case. This requirement severely limits the scalability of classical explainability methods, as traditional techniques such as sampling and model-agnostic surrogates introduce unavoidable trade-offs between computational efficiency and explanation fidelity. To address the inherent bottlenecks of traditional computation, researchers have increasingly turned to quantum computing as an alternative paradigm (Brassard et al., 2002; Montanaro, 2016). Recent work (Burge et al., 2025) demonstrates a quantum algorithm that encodes coalition weights and marginal contributions into quantum states while employing amplitude amplification to achieve quadratic speedup relative to classical Monte Carlo strategies. By adapting this algorithm to the graph domain, it becomes feasible to compute Shapley values at the subgraph level with all possible coalitions of size $2^n$. This advancement effectively mitigates the exponential coalition bottleneck, thereby enhancing both the scalability and explainability of Shapley-based explanations in graph neural networks.

## 3 Background

### 3.1 Graph Neural Networks (GNNs)

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E}, X)$ denote an undirected graph with node set $\mathcal{V}$, edge set $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$, and node feature matrix $X \in \mathbb{R}^{n \times d}$, where $n$ is the number of nodes and $d$ is the dimensionality of node features. A GNN (Gilmer

et al., 2017) computes node embeddings by iteratively aggregating neighborhood information:

$$\mathbf{h}_v^{(l)} = U^{(l)}\left(\mathbf{h}_v^{(l-1)}, \sum_{u \in \mathcal{N}(v)} M^{(l)}\left(\mathbf{h}_v^{(l-1)}, \mathbf{h}_u^{(l-1)}\right)\right) \tag{1}$$

where, $M^{(l)}$ and $U^{(l)}$ are (possibly learnable) functions, and $\mathbf{h}_v^{(0)} = \mathbf{x}_v$. For graph-level tasks, node embeddings are aggregated via a readout function $R$: $\mathbf{h}_{\mathcal{G}} = R(\{\!\{\mathbf{h}_v^{(L)}\}\!\})$.

## 3.2 Shapley Values for Explanations

A *coalitional game* (Winter, 2002) is given by $(P, v)$, where $P = \{1, \ldots, n\}$ and a function, $v : 2^P \to \mathbb{R}$ assigns a value to each coalition $S \subseteq P$, with $v(\emptyset) = 0$. For a player $p_j \in P$, the Shapley value is defined as

$$\phi(p_j) = \sum_{S \subseteq P \setminus \{p_j\}} w(|S|, n)\left[v(S \cup \{p_j\}) - v(S)\right] \tag{2}$$

where $w(|S|, n) = \frac{1}{\binom{n-1}{|S|}} \cdot \frac{1}{n}$ is the coalition weighting term. The quantity $\phi(p_j)$ is the unique solution that satisfies the following axioms:

- **Efficiency:** The total value is distributed, i.e., $\sum_{p_j \in P} \phi(p_j) = v(P)$.

- **Symmetry:** If two players $p_j, p_k$ satisfy $v(S \cup \{p_j\}) = v(S \cup \{p_k\})$ for all $S \subseteq P \setminus \{p_j, p_k\}$, then $\phi(p_j) = \phi(p_k)$.

- **Dummy:** If a player $p_j$ satisfies $v(S \cup \{p_j\}) = v(S)$ for all $S \subseteq P \setminus \{p_j\}$, then $\phi(p_j) = 0$.

- **Additivity:** For games with value functions $v$ and $v'$, the Shapley value for the summed game is $\phi(p_j)(v + v') = \phi(p_j)(v) + \phi(p_j)(v')$.

## 3.3 Quantum Estimation of the Shapley Value

Consider the classical coalitional game $(P, v)$ with $n$ players. In the quantum algorithm setup (Burge et al., 2025), represent the game as $(P, U)$, where $U : 2^P \to [0, 1]$ is a normalized utility function that maps each coalition $S \subseteq P$ to a value in $[0, 1]$. The goal is to approximate the Shapley value $\phi(p_j)$ of participant $p_j \in P$ with additive error $\varepsilon$. The algorithm employs three quantum registers: *Partition register* $Q_{\mathrm{pt}}$ with $\ell$ qubits, used to encode an amplitude distribution proportional to the coalition-weight coefficients, $\omega(n, r)$, where $r = |S|$ and:

$$\ell = \mathcal{O}\left(\log \frac{(U_{\max} - U_{\min})\, n}{\varepsilon}\right); \tag{3}$$

The corresponding bounds for the normalized utility function $U$ are defined as

$$U_{\max} = \max_{S \subseteq P} U(S), \qquad U_{\min} = \min_{S \subseteq P} U(S). \tag{4}$$

*Player register* $Q_{\mathrm{pl}}$ with $n$ qubits, stores a superposition of all coalitions $S \subseteq P \setminus \{p_j\}$; and, the *Utility register* $Q_{\mathrm{ut}}$ with a single qubit, represents the normalized utility of each coalition. Controlled rotations, $R_j$, parameterized by the partition amplitudes, prepare the superposition

$$\sum_{S \subseteq P \setminus \{p_j\}} \sqrt{\omega(n, |S|)}\, |S\rangle_{Q_{\mathrm{pl}}}, \tag{5}$$

so that, the amplitude of each coalition corresponds to its Shapley weight. Two quantum oracles, $U_{\mathrm{val}}^{(+)}$ and $U_{\mathrm{val}}^{(-)}$, implement the normalized utility function $U$, conditioning on whether $p_j$ is included (+) or excluded

($-$) from the coalition. Applying the quantum amplitude estimation routine described in Montanaro (2015) to these states yields the quantities $\phi^{(+)}(p_j)$ and $\phi^{(-)}(p_j)$, which correspond to the expected marginal contributions of $p_j$, when it is included in, and excluded from a coalition, respectively. The Shapley value is then obtained as

$$\phi(p_j) = \phi^{(+)}(p_j) - \phi^{(-)}(p_j) \tag{6}$$

with total error bounded by $\varepsilon$ and overall query complexity $\mathcal{O}\left(\frac{U_{\max} - U_{\min}}{\varepsilon}\right)$, providing a near-quadratic speedup compared to classical Monte Carlo estimation. The speedup is achieved because quantum amplitude estimation requires only $\mathcal{O}(1/\epsilon)$ queries to reach an additive error $\epsilon$, compared to $\mathcal{O}(1/\epsilon^2)$ queries for classical Monte Carlo approaches, thus substantially reducing computational costs. Formal proofs of correctness, error bounds, and complexity guarantees for the proposed quantum procedures are presented in Montanaro (2015); Burge et al. (2025).
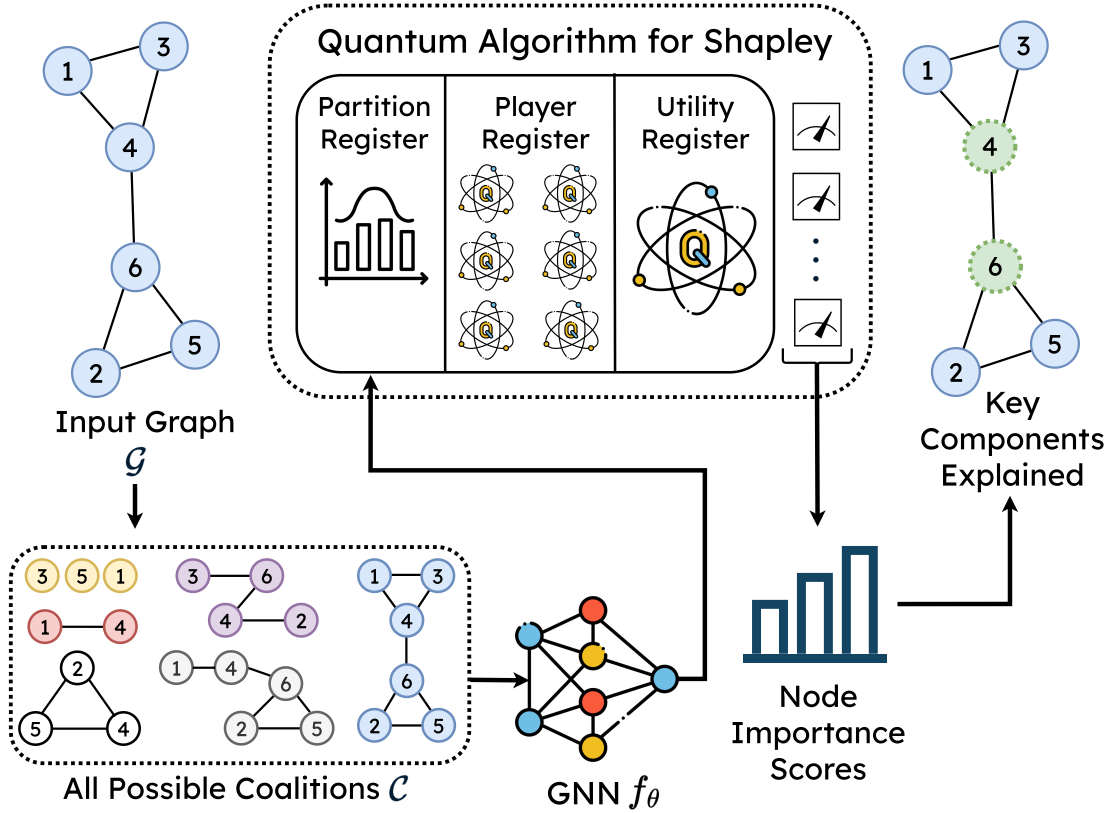
## 4  QGShap



Figure 1: QGSHAP Workflow: The input graph $\mathcal{G}$ is mapped to node coalitions $\mathcal{C}$, each scored by the trained GNN $f_\theta$ via a masking oracle to obtain coalition values. A quantum module prepares three registers: a *Partition* register encoding Shapley weights, a *Player* register encoding coalition indices, and a *Utility* register encoding normalized coalition scores. Quantum amplitude estimation over the *Utility* register aggregates weighted marginal contributions, yielding node-level Shapley attributions as the final explanations.

We present QGSHAP, a post hoc explainability framework that leverages quantum speedup to obtain exact Shapley value explanations for GNN predictions. Starting from a trained GNN, $f_\theta$, and an input graph, $G = (V, E)$, we exhaustively enumerate all non-empty node subsets, $S \subseteq V$, and generate corresponding masked graphs, $G_S$, via zero-fill encoding, where excluded nodes are replaced with zero vectors. Each masked graph is evaluated by $f_\theta$ to compute the cooperative game-theoretic value $v(S)$, ensuring that every node's marginal contribution is precisely captured without resorting to sampling heuristics.
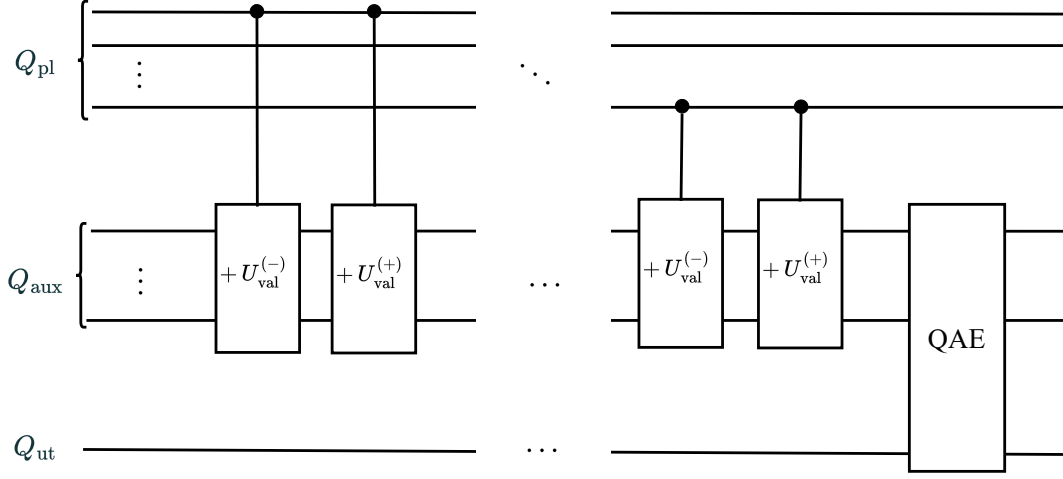
Figure 2: Circuit of the QGSHAP utility oracle $U_{\text{val}}^{(\pm)}$. The player register $Q_{\text{pl}}$ encodes coalitions $S \subseteq V \setminus \{p_j\}$, while the auxiliary register $Q_{\text{aux}}$ and utility register $Q_{\text{ut}}$ store normalized cooperative values $\hat{v}(S)$ and $\hat{v}(S \cup \{p_j\})$. The quantum oracles $U_{\text{val}}^{(-)}$ and $U_{\text{val}}^{(+)}$ correspond to evaluating coalitions without and with node $p_j$, respectively. Quantum Amplitude Estimation (QAE) is then applied to $Q_{\text{ut}}$ to obtain the expected contributions $\phi^{(+)}(p_j)$ and $\phi^{(-)}(p_j)$, which are combined to reconstruct the exact Shapley value $\phi(p_j)$ as described in Section 3.3.

Building on the quantum Shapley value estimation framework introduced in Section 3.3 and originally proposed in (Burge et al., 2025) (Section 5), we adopt the same state-preparation procedure for encoding exact Shapley weights. For the amplitude estimation step, we employ the quantum routine described in (Montanaro, 2015). Accordingly, to prepare the exact Shapley weights in a quantum state, we first normalize the cooperative values:

$$\hat{v}(S) = \frac{v(S) - \min_{S'} v(S')}{\max_{S'} v(S') - \min_{S'} v(S')} \in [0, 1]. \tag{7}$$

We then allocate a *player register* $Q_{\text{pl}}$ of $|V|$ qubits, encoding coalition membership, a *partition register* $Q_{\text{pt}}$ initialized via beta-function rotations to load amplitudes proportional to the Shapley coefficients $w_{|S|,|V|}$, and a *utility register* $Q_{\text{ut}}$ to store the normalized cooperative value $\hat{v}(S)$ for each coalition. Controlled rotations between the partition and player registers prepare a superposition in which the amplitude of each basis state $|S\rangle$ is proportional to $\sqrt{w_{|S|,|V|}}$. We then invoke the quantum amplitude estimation subroutine to extract the weighted expected contributions $\phi^{(+)}(p_j)$ and $\phi^{(-)}(p_j)$ from the utility register for each participant node $p_j$, achieving a quadratic reduction in sampling complexity compared to classical Monte Carlo methods. Finally, we reconstruct the Shapley value of each node by computing the difference in weighted expected values and denormalizing:

$$\phi(p_j) = \left( \max_S v(S) - \min_S v(S) \right) \left( \phi^{(+)}(p_j) - \phi^{(-)}(p_j) \right). \tag{8}$$

Then, we normalize again across all nodes to produce a global importance ranking. Thus, QGSHAP computes Shapley contributions per node within each coalition, treating nodes as "players" in many small cooperative games, rather than treating entire coalitions or subgraphs as atomic units. For every coalition generated via exhaustive enumeration, it constructs all subsets of nodes within that coalition, evaluates the model's value

function for each subset, and applies quantum amplitude estimation to compute the exact marginal Shapley value for each individual node in that coalition. This hierarchical, node-centric approach differs fundamentally from SubgraphX, which treats sampled subgraphs and residual nodes as players, simultaneously, as well as, GraphSHAP-IQ, which derives exact any-order Shapley interactions across nodes within receptive fields, but without explicit coalition traversal. The quantum speedup from amplitude estimation (Montanaro, 2015) over normalized subset values within coalitions yields unbiased, high-fidelity node attributions in $\mathcal{O}(1/\epsilon)$ complexity, surpassing classical Monte Carlo's $\mathcal{O}(1/\epsilon^2)$, and motivates its application for precise quantum GNN explainability by capturing subtle intra-coalition interactions overlooked by SubgraphX's subgraph scoring and GraphSHAP-IQ's receptive-field constraints.

QGSHAP's pipeline is illustrated in Fig. 1, the circuit of the quantum routine in Fig. 2, and we formally detail QGSHAP in Algorithm 1.

---

**Algorithm 1** QGSHAP: Exact Shapley Value Estimation

---

1: **Input:** Trained GNN $f_\theta$; input graph $G = (V, E)$
2: **Output:** Node-level Shapley value explanations $\{\phi_i\}_{i \in V}$
3: Enumerate all non-empty subsets $\mathcal{C} = \{S \subseteq V : S \neq \emptyset\}$
4: **for all** $S \in \mathcal{C}$ **do**
5:     Construct masked graph $G_S$ via zero-fill encoding
6:     Evaluate cooperative value $v(S) = f_\theta(G_S)$
7: **end for**
8: Compute $v_{\min} = \min_{S \in \mathcal{C}} v(S)$ and $v_{\max} = \max_{S \in \mathcal{C}} v(S)$
9: Normalize all cooperative values: $\hat{v}(S) = \frac{v(S) - v_{\min}}{v_{\max} - v_{\min}}$ for all $S \in \mathcal{C}$
10: Allocate player register $Q_{\text{pl}}$ ($|V|$ qubits), partition register $Q_{\text{pt}}$, and utility register $Q_{\text{ut}}$
11: Encode Shapley weights $w_{|S|,|V|}$ via beta-rotation circuits in $Q_{\text{pt}}$
12: Apply controlled rotations $R_j$ to prepare superposition of coalitions in $Q_{\text{pl}}$ and store $\hat{v}(S)$ in $Q_{\text{ut}}$
13: **for all** $p_j \in V$ **do**
14:     Apply quantum amplitude estimation on $Q_{\text{ut}}$ to obtain $\phi^{(+)}(p_j), \phi^{(-)}(p_j)$
15:     $\phi(p_j) \leftarrow (v_{\max} - v_{\min}) \cdot \left(\phi^{(+)}(p_j) - \phi^{(-)}(p_j)\right)$
16: **end for**
17: Normalize $\{\phi(p_j)\}$ to produce final node importance scores

---

## 5 Experiments

To thoroughly evaluate our proposed approach, we construct a controlled experimental setup that supports both quantitative and qualitative analysis of explanation performance. This section outlines the datasets, details the model used for prediction tasks, and describes the explanation methodology along with evaluation metrics applied to assess explanation effectiveness. Our implementation is available at https://github.com/smlab-niser/qgshap.

### 5.1 Datasets

**Bridge:** The Bridge detection dataset (Toyokuni & Yamada, 2023) consists of synthetically generated graphs formed by connecting two cycle graphs (3-5 nodes each) via a bridge linking selected nodes. Node identities are randomized across samples to ensure the model and explanations rely on graph structure rather than fixed node positions. The training set includes configurations with up to 15 nodes, containing both graphs with a bridge edge (label 1) and disconnected cycle pairs (label 0), totaling 60 balanced samples. The test set comprises 20 graphs with bridge edges, with exactly 8 nodes across four fixed configurations: (3+3), (3+4), (4+3), and (4+4), with each number denoting the order of the cycles connected with a bridge, enabling evaluation of generalization to unseen structures.

**BA2-Motif:** The BA2-Motif dataset (Ying et al., 2019) extends the Barabási-Albert (BA) model into which exactly one motif is inserted either a house (label 1) or a cycle (label 0). This dataset introduces motif-level

explanation within scale-free topologies. Following the *ExplainerDataset* setup in PyTorch Geometric, 50 train and 50 test graphs are generated. Notably, the house motif differs from a 5-cycle graph by including one additional edge closing the structure, designated as the 'house edge'. The ability of an explainer to accurately pinpoint and explain the house edge becomes an explicit test for evaluating motif-level ground truth recovery in graph explanations.

## 5.2 Model Training

We implement a Graph Isomorphism Network (GIN) (Xu et al., 2019) classifier using PyTorch 2.8.0[2] and PyTorch Geometric 2.6.1[3] (CUDA 12.8). The model consists of a multi-layer perceptron (MLP) as an encoder layer, three GIN layers, followed by another MLP as a decoder. We set the hidden dimension to 128. Training is conducted for 100 epochs using the Adam optimizer (Kingma & Ba, 2015) with a learning rate of $10^{-3}$. The model is optimized using binary cross-entropy loss between the predicted and true class labels. All test graphs in this study were limited to 8 nodes to ensure feasibility. In the quantum setting, each node requires a qubit in the player register, and the partition register must also scale with the number of nodes to encode coalition-weight amplitudes. Beyond 8 nodes, the number of qubits and circuit complexity grow rapidly, making exact Shapley value estimation infeasible with current quantum simulation resources. We employ the publicly available reference implementation[4] to compute quantum-accelerated Shapley values for our experiments, executing all quantum subroutines on a Qiskit-based simulator[5].

## 5.3 Evaluation Metrics

We evaluate the quality of explanations using a set of standard metrics commonly used in the GNN explainability literature (Ying et al., 2019), each capturing different aspects of explanatory performance.

**Top-$k$ Accuracy** This metric measures the frequency with which the ground-truth target nodes appear among the top-$k$ most important nodes, as ranked by the explanation model. Formally, it is defined as:

$$\text{Acc}_{\text{top}-k} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}\left[\text{target}_i \in \text{Top-}k(\phi(p_j))\right],$$

where, $\phi(p_j)$ represents the importance scores assigned to nodes in instance $p_j$ (Shapley values for Shapley-based explainers, or standard node importance scores otherwise), and $\mathbb{I}[\cdot]$ is the indicator function. Only top-2 accuracy ($k = 2$) is reported, since for both Bridge and BA2-Motif datasets, we focus on whether the explainer identifies either the bridge edge or the 'house edge' that completes the house motif from a 5-cycle. Higher values indicate better alignment between the explainer's output and the true important nodes.

**Fidelity** Fidelity (Ying et al., 2019) evaluates how well the explanation aligns with the model's predictions when key nodes are selectively retained or removed.

- *Fidelity-plus* ($\text{Fid}^+$) measures the model's confidence or prediction consistency, when only the top-$k$ important nodes (by $\phi_i$) are retained, while all others are removed.

- *Fidelity-minus* ($\text{Fid}^-$) measures the effect of removing the top-$k$ nodes while keeping the rest.

Formally, if $G$ is the original graph, $S$ is the set of top-$k$ important nodes, and $y_c$ is the predicted class on $G$, then $\text{Fid}^+$ is defined as

$$\text{Fid}^+ = P_{\text{keep}}(y_c) - P_{\text{base}}(y_c),$$

where $P_{\text{base}}(y_c)$ is the model's predicted class probability on the full graph $G$, and $P_{\text{keep}}(y_c)$ is the probability on the induced subgraph containing only $S$. Conversely,

$$\text{Fid}^- = P_{\text{base}}(y_c) - P_{\text{remove}}(y_c),$$

---

[2] https://pytorch.org/
[3] https://pytorch-geometric.readthedocs.io/en/2.6.1/
[4] https://github.com/iain-burge/QuantumShapleyValueAlgorithm
[5] https://www.ibm.com/quantum/qiskit

where $P_{\text{remove}}(y_c)$ is the predicted probability for class $y_c$ on the complement graph $G \setminus S$.

An effective explainer should yield high $\mathsf{Fid}^+$ (the top-ranked nodes alone suffice to reproduce the model's prediction) and low $\mathsf{Fid}^-$ (removing these nodes substantially changes the prediction), reflecting high confidence that it has correctly identified the most influential nodes driving the model's decision.

**Sparsity**   Sparsity (Ying et al., 2019) captures how concise the explanation is by computing the proportion of nodes that receive low importance scores:

$$S = 1 - \frac{|\{p_j : \phi(p_j) \geq 0.1 \max_i \phi(p_i)\}|}{N}.$$

Here, $\phi(p_j)$ denotes the importance score of node $p_j$ (Shapley value for Shapley-based explainers, or the explainer's node importance score otherwise), and $N$ is the total number of nodes. Higher sparsity indicates that the explanation focuses on a small set of high-importance nodes, thereby improving explainability and reducing noise.

**Graph Explanation Accuracy (GEA)**   GEA (Agarwal et al., 2023) quantifies how closely an explainer's predicted important nodes match the ground truth nodes in a graph, using the Jaccard similarity index. It is defined as:

$$\text{GEA} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}}$$

where, TP (true positives) is the number of nodes correctly identified as important, FP (false positives) is the number of nodes wrongly identified as important, and FN (false negatives) is the number of true important nodes missed. GEA yields a value between 0 (no overlap) and 1 (perfect match), offering an intuitive, symmetric measure of explanatory set quality by considering both types of errors equally. Moreover, unlike $\mathsf{Fid}^+$ or $\mathsf{Fid}^-$, which can remain high even when an explainer assigns maximal score to the wrong nodes due to distributional sufficiency effects, GEA directly penalizes such target misalignments by measuring set overlap with ground truth, thereby detecting cases where high fidelity coexists with incorrect node attributions.

### 5.4   Results and Discussion

Visualizations for the Bridge and BA2-Motif datasets are provided in Fig. 3 and 4, with quantitative metrics summarized in Table 1. We show graph heatmaps only for SubGraphX, as it is the sole competitive Shapley-based explainer among the compared methods.

Table 1: Comparison of explanation metrics for different explainers.

| Dataset | Explainer | $\mathsf{Fid}^+$ | $\mathsf{Fid}^-$ | Sparsity | GEA | Top-2 Acc |
|---------|-----------|------------------|------------------|----------|-----|-----------|
| Bridge | GNNExplainer | $0.60 \pm 0.49$ | $1.00 \pm 0.00$ | $0.75 \pm 0.00$ | $\underline{0.07 \pm 0.13}$ | $\underline{0.10 \pm 0.20}$ |
| | PGExplainer | $0.99 \pm 0.00$ | $1.00 \pm 0.00$ | $0.75 \pm 0.00$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | SubgraphX | $\mathbf{1.00 \pm 0.00}$ | $1.00 \pm 0.00$ | $0.75 \pm 0.00$ | $\mathbf{1.00 \pm 0.00}$ | $\mathbf{1.00 \pm 0.00}$ |
| | **QGShap (Ours)** | $\mathbf{1.00 \pm 0.00}$ | $1.00 \pm 0.00$ | $0.75 \pm 0.00$ | $\mathbf{1.00 \pm 0.00}$ | $\mathbf{1.00 \pm 0.00}$ |
| BA2-Motif | GNNExplainer | $\mathbf{1.00 \pm 0.00}$ | $1.00 \pm 0.00$ | $0.75 \pm 0.00$ | $\underline{0.32 \pm 0.11}$ | $0.50 \pm 0.19$ |
| | PGExplainer | $0.01 \pm 0.01$ | $1.00 \pm 0.00$ | $0.75 \pm 0.00$ | $0.00 \pm 0.00$ | $0.00 \pm 0.00$ |
| | SubgraphX | $\mathbf{1.00 \pm 0.00}$ | $1.00 \pm 0.00$ | $0.75 \pm 0.00$ | $\mathbf{0.40 \pm 0.00}$ | $\underline{0.79 \pm 0.25}$ |
| | **QGShap (Ours)** | $\mathbf{1.00 \pm 0.00}$ | $1.00 \pm 0.00$ | $0.75 \pm 0.00$ | $\mathbf{0.40 \pm 0.00}$ | $\mathbf{1.00 \pm 0.00}$ |

**Bridge**   On the Bridge dataset, our proposed QGSHAP consistently achieved perfect performance across all evaluation metrics, reflecting both precision and reliability in its explanations. As shown in Table 1, QGSHAP attained Fidelity$^+$, GEA, and Top-2 Accuracy scores of $1.00 \pm 0.00$, matching or surpassing the strongest baseline, SubgraphX. In every case, it successfully identified the critical bridge nodes driving the model's predictions, demonstrating clear explainability and stability. These results highlight QGSHAP's ability to deliver faithful and consistent explanations that align closely with the underlying graph structure and decision logic.
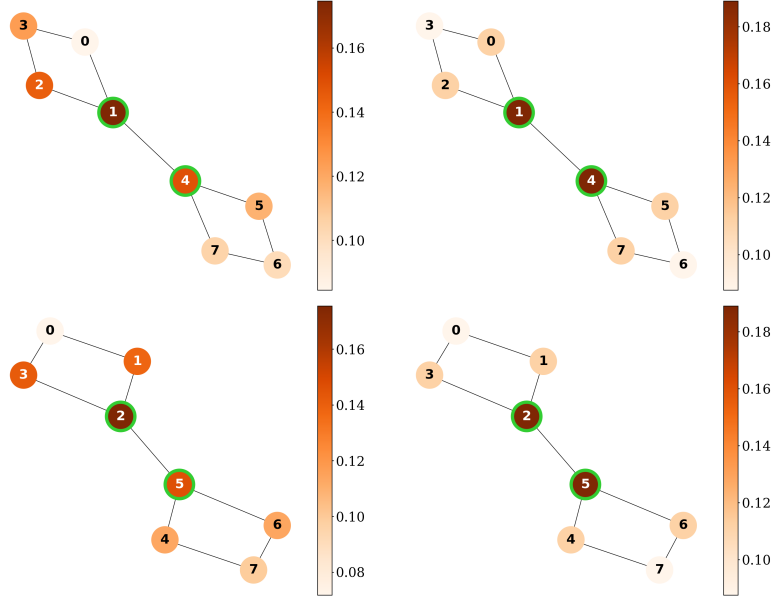
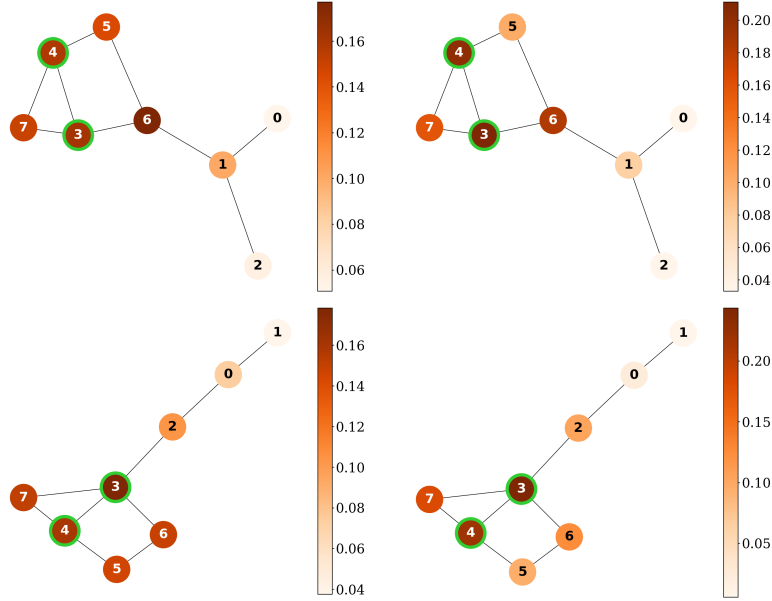Figure 3: **Bridge**: Subgraph explanations using SubgraphX and QGSʜᴀᴘ. **Left -** SubgraphX . **Right -** QGSʜᴀᴘ



Figure 4: **BA2-Motif**: Subgraph explanations using SubgraphX and QGSʜᴀᴘ. **Left -** SubgraphX . **Right -** QGSʜᴀᴘ

**BA2-Motif** For the BA2-Motif dataset, QGSʜᴀᴘ also produced highly competitive and insightful results, achieving the highest Top-2 Accuracy $(1.00 \pm 0.00)$ among all explainers. In most instances, it correctly pinpointed the two key nodes and their connecting edge responsible for predicting the 'house' motif, capturing the core structural reasoning of the model. Even in a few challenging cases, where not all motif nodes were ranked at the top, QGSʜᴀᴘ consistently prioritized the most influential node pairs linked to the correct class. This behavior underscores its robustness and interpretive strength in highlighting the decisive substructures within complex graph motifs.

## 6 Conclusion

We introduce QGSHAP, a post hoc explainability framework that combines cooperative game theory and quantum computation to produce exact Shapley value explanations for GNN predictions. Unlike classical sampling or approximation-based methods, QGSHAP evaluates all coalitions, capturing node influence precisely and verifiably rather than relying on heuristics. Although currently limited to small graphs by hardware constraints, QGSHAP demonstrates that quantum computation can make exact, *classically intractable* Shapley calculations practical and establishes a benchmark for evaluating classical explainers while bridging explainable GNNs with developments in quantum computing. Extending both the GNN and explanation modules into the quantum domain, our framework provides a principled and scalable approach to explainability in GNNs, classical or quantum. Future work could explore iterative, noise-resilient amplitude amplification strategies for robustness under realistic hardware constraints and fault-tolerant settings.

### Limitations

Although the method offers a near-quadratic speedup over classical Monte Carlo techniques up to polylogarithmic factors, QGSHAP remains constrained to small graphs due to the exponential number of coalitions and gate preparations, which increases qubit requirements and circuit depth. The practical cost of classical simulation remains high: even on a system with a 48-core AMD CPU and an Nvidia A100 GPU (80 GiB VRAM), the Bridge and BA2-Motif experiments required approximately 31 and 42 hours of runtime, respectively. Moreover, quantum noise and decoherence can reduce the precision of amplitude estimation, and access to high-quality quantum hardware remains limited. These factors underscore the current operational scope of QGSHAP and highlight the need for advances that enable scaling to larger and more complex graphs.

## References

Chirag Agarwal, Owen Queen, Himabindu Lakkaraju, and Marinka Zitnik. Evaluating explainability for graph neural networks. *Scientific Data*, 10(1):144, 2023.

Selahattin Akkas and Ariful Azad. Gnnshap: Scalable and accurate gnn explanation using shapley values. In *Proceedings of the ACM Web Conference 2024*, pp. 827–838. ACM, May 2024. doi: 10.1145/3589334.3645599.

Meghyn Bienvenu, Diego Figueira, and Pierre Lafourcade. When is shapley value computation a matter of counting? *Proc. ACM Manag. Data*, 2(2), May 2024. doi: 10.1145/3651606.

Gilles Brassard, Peter Høyer, Michele Mosca, and Alain Tapp. Quantum amplitude amplification and estimation, 2002. ISSN 0271-4132.

Iain Burge, Michel Barbeau, and Joaquin Garcia-Alfaro. A shapley value estimation speedup for efficient explainable quantum ai. *Preprint*, April 2025. doi: 10.48550/arXiv.2412.14639.

Xiaotie Deng and Christos H. Papadimitriou. On the complexity of cooperative solution concepts. *Math. Oper. Res.*, 19(2):257–266, May 1994. ISSN 0364-765X. doi: 10.1287/moor.19.2.257.

Alexandre Duval and Fragkiskos D Malliaros. Graphsvx: Shapley value explanations for graph neural networks. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 302–318. Springer, 2021.

Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pp. 1263–1272. Pmlr, 2017.

Qiang Huang, Makoto Yamada, Yuan Tian, Dinesh Singh, and Yi Chang. Graphlime: Local interpretable model explanations for graph neural networks. *IEEE Transactions on Knowledge and Data Engineering*, 35(7):6968–6972, 2022.

Rucha Bhalchandra Joshi and Subhankar Mishra. Learning graph representations. In Anupam Biswas, Ripon Patgiri, and Bhaskar Biswas (eds.), *Principles of Social Networking: The New Horizon and Emerging Challenges*, pp. 209–228. Springer Singapore, Singapore, 2022. ISBN 978-981-16-3398-0. doi: 10.1007/978-981-16-3398-0_10.

Majd Khalil and Benny Kimelfeld. The Complexity of the Shapley Value for Regular Path Queries. In Floris Geerts and Brecht Vandevoort (eds.), *26th International Conference on Database Theory (ICDT 2023)*, volume 255 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 11:1–11:19, Dagstuhl, Germany, 2023. Schloss Dagstuhl – Leibniz-Zentrum für Informatik. ISBN 978-3-95977-270-9. doi: 10. 4230/LIPIcs.ICDT.2023.11.

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. URL https://arxiv.org/abs/1412.6980.

Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL https://openreview.net/forum?id=SJU4ayYgl.

H. W. Kuhn and A. W. Tucker. *Contributions to the Theory of Games, Volume II*, volume 28 of *Annals of Mathematics Studies*. Princeton University Press, 1953.

Xiao Li, Li Sun, Mengjie Ling, and Yan Peng. A survey of graph neural network based recommendation in social networks. *Neurocomputing*, 549:126441, 2023.

Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, pp. 4768–4777, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.

Dongsheng Luo, Wei Cheng, Dongkuan Xu, Wenchao Yu, Bo Zong, Haifeng Chen, and Xiang Zhang. Parameterized explainer for graph neural network. *Advances in neural information processing systems*, 33: 19620–19631, 2020.

Peter Mernyei, Konstantinos Meichanetzidis, and Ismail Ilkan Ceylan. Equivariant quantum graph circuits. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 15401–15420. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/mernyei22a.html.

Ashley Montanaro. Quantum speedup of monte carlo methods. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 471(2181):20150301, September 2015. ISSN 1471-2946. doi: 10.1098/rspa.2015.0301.

Ashley Montanaro. Quantum algorithms: an overview. *npj Quantum Information*, 2(1), January 2016. ISSN 2056-6387. doi: 10.1038/npjqi.2015.23. URL http://dx.doi.org/10.1038/npjqi.2015.23.

Maximilian Muschalik, Fabian Fumagalli, Paolo Frazzetto, Janine Strotherm, Luca Hermes, Alessandro Sperduti, Eyke Hüllermeier, and Barbara Hammer. Exact computation of any-order shapley interactions for graph neural networks. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=9tKC0YM8sX.

Phillip E Pope, Soheil Kolouri, Mohammad Rostami, Charles E Martin, and Heiko Hoffmann. Explainability methods for graph convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 10772–10781, 2019.

Patrick Reiser, Marlen Neubert, André Eberhard, Luca Torresi, Chen Zhou, Chen Shao, Houssam Metni, Clint van Hoesel, Henrik Schopmans, Timo Sommer, et al. Graph neural networks for materials science and chemistry. *Communications Materials*, 3(1):93, 2022.

Robert Schwarzenberg, Marc Hübner, David Harbecke, Christoph Alt, and Leonhard Hennig. Layerwise relevance visualization in convolutional text graph classifiers. In Dmitry Ustalov, Swapna Somasundaran, Peter Jansen, Goran Glavaš, Martin Riedl, Mihai Surdeanu, and Michalis Vazirgiannis (eds.), *Proceedings of the Thirteenth Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-13)*, pp. 58–62, Hong Kong, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-5308. URL https://aclanthology.org/D19-5308/.

David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy P. Lillicrap, Fan Hui, L. Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of go without human knowledge. *Nature*, 550:354–359, 2017. URL https://api.semanticscholar.org/CorpusID:205261034.

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *International Conference on Learning Representations*, 2015.

Ayato Toyokuni and Makoto Yamada. Structural explanations for graph neural networks using hsic. *arXiv preprint arXiv:2302.02139*, 2023.

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=rJXMpikCZ.

Eyal Winter. The Shapley value. In *Handbook of Game Theory with Economic Applications*, volume 3, pp. 2025–2054. Elsevier, 2002. doi: 10.1016/S1574-0005(02)03016-3.

Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph neural networks in recommender systems: a survey. *ACM Computing Surveys*, 55(5):1–37, 2022.

Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1): 4–24, January 2021. ISSN 2162-2388. doi: 10.1109/tnnls.2020.2978386.

Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=ryGs6iA5Km.

Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gnnexplainer: Generating explanations for graph neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

Hao Yuan, Haiyang Yu, Jie Wang, Kang Li, and Shuiwang Ji. On explainability of graph neural networks via subgraph explorations. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 12241–12252. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/yuan21c.html.

Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pp. 818–833. Springer, 2014.

Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2921–2929, 2016.