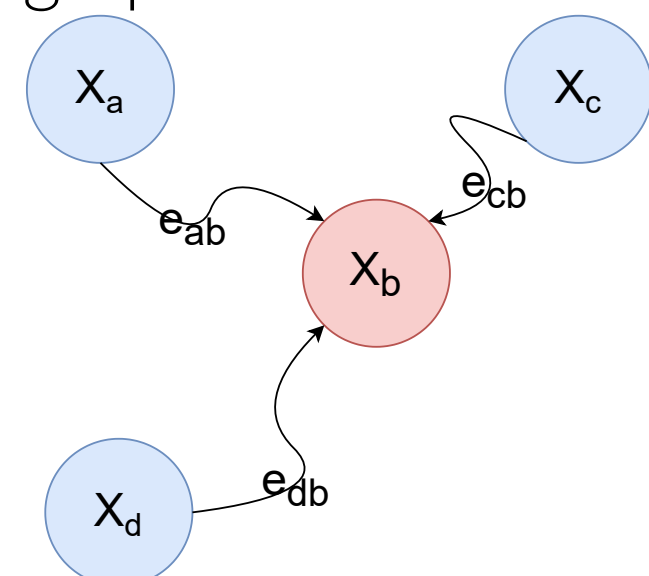


Background

Graph Neural Networks (GNNs) are Neural Networks for graph data:

$$h_b^0 = x_b \in \mathbb{R}^d$$

$$h_b^l = \text{Upd}^l(h_b^{l-1}, \text{Aggr}^l(\{h_u^{l-1} : u \in N(b)\}))$$



Where Upd is a Neural Network and Aggr is any permutation invariant function.

Motivation

GNNs lack interpretability, thus hindering understanding, debugging, and human trust:

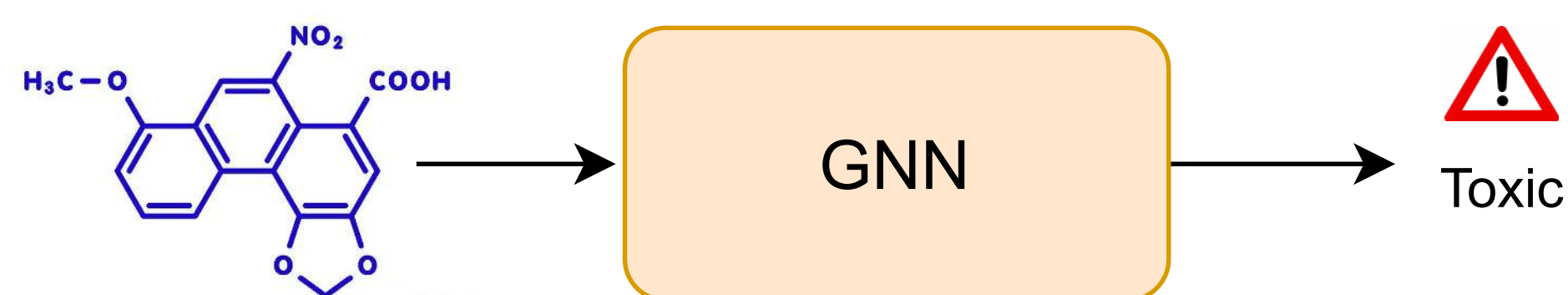


Figure 1. Aristolochic acids are a family of carcinogenic, mutagenic, and nephrotoxic phytochemicals commonly found in the flowering plant family Aristolochiaceae.

As popular post-hoc methods have been found to fall short in reliably explaining trained GNNs [5, 4], new **explainable by-design** architectures have been proposed:

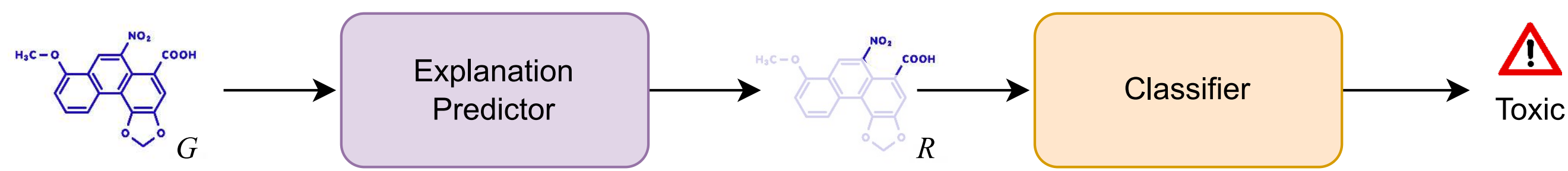


Figure 2. Pipeline of Self-Explainable GNNs (SEGNNs).

⚠ Nonetheless, some SEGNNs are found to be more **faithful** to random explanations than to their true explanations [3]. ⚠

Our contribution

We aim to study the root of this issue while providing insights into how to build more reliable SEGNNs:

- **RQ1:** What characterizes a good explanation?
- **RQ2:** How good are SEGNNs?
- **RQ3:** Can we go beyond subgraph-based explanations?

RQ1: What characterizes a good explanation? [1]

Current literature measures *how much the model adheres to its explanation* by measuring the **faithfulness** of explanations:

- *sufficient*, i.e., keeping it fixed shields the model's output from changes to its complement $C = G \setminus R$

$$SUF_{d,p_R}(R) = \mathbb{E}_{G' \sim p_R}[\Delta_d(G, G')],$$

- *necessary*, i.e., altering it affects the model's output even with C fixed

$$NEC_{d,p_C}(R) = \mathbb{E}_{G' \sim p_C}[\Delta_d(G, G')]$$

ⓘ We provide a taxonomy of the current faithfulness metric:

Table 1. SUF and NEC recover existing faithfulness metrics for appropriate choices of divergence d and interventional distributions p_R and p_C .

Metric	Estimates	Divergence d	Allowed changes
Unf	Suf	$KL(p_\theta(\cdot G), p_\theta(\cdot G'))$	zero out all irrelevant features
Fid-	RFid-	$ p_\theta(\hat{y} G) - p_\theta(\hat{y} G') $	zero out all irrelevant features, delete all irrelevant edges
PS		$\mathbb{1}\{p_\theta(\hat{y} G) = p_\theta(\hat{y} G')\}$	delete a random subset of irrelevant edges
Fid+	Nec	$ p_\theta(\hat{y} G) - p_\theta(\hat{y} G') $	multiply all irrelevant elements by relevance scores
RFid+			zero out all relevant features, delete all relevant edges
PN		$\mathbb{1}\{p_\theta(\hat{y} G) \neq p_\theta(\hat{y} G')\}$	delete a random subset of relevant edges
			multiply all relevant elements by relevance scores

Table 2. Model ranking and SUF results according to different p_R .

Split	Model	Motif2	
		$p_R^{id_1}$	$p_R^{id_2}$
	LECI	1 (81 ± 03)	2 (82 ± 03)
ID	GSAT	2 (78 ± 01)	1 (84 ± 02)
	CIGA	3 (65 ± 07)	3 (73 ± 06)

ⓘ Metrics are not interchangeable in the sense that metric values across different metric parameters are not consistent.

RQ1: What characterizes a good explanation? (cont.) [1]

- ⓘ Previous Necessity metrics do not penalize useless explanations
- ⓘ We propose a new necessity metric that penalizes overly large explanations

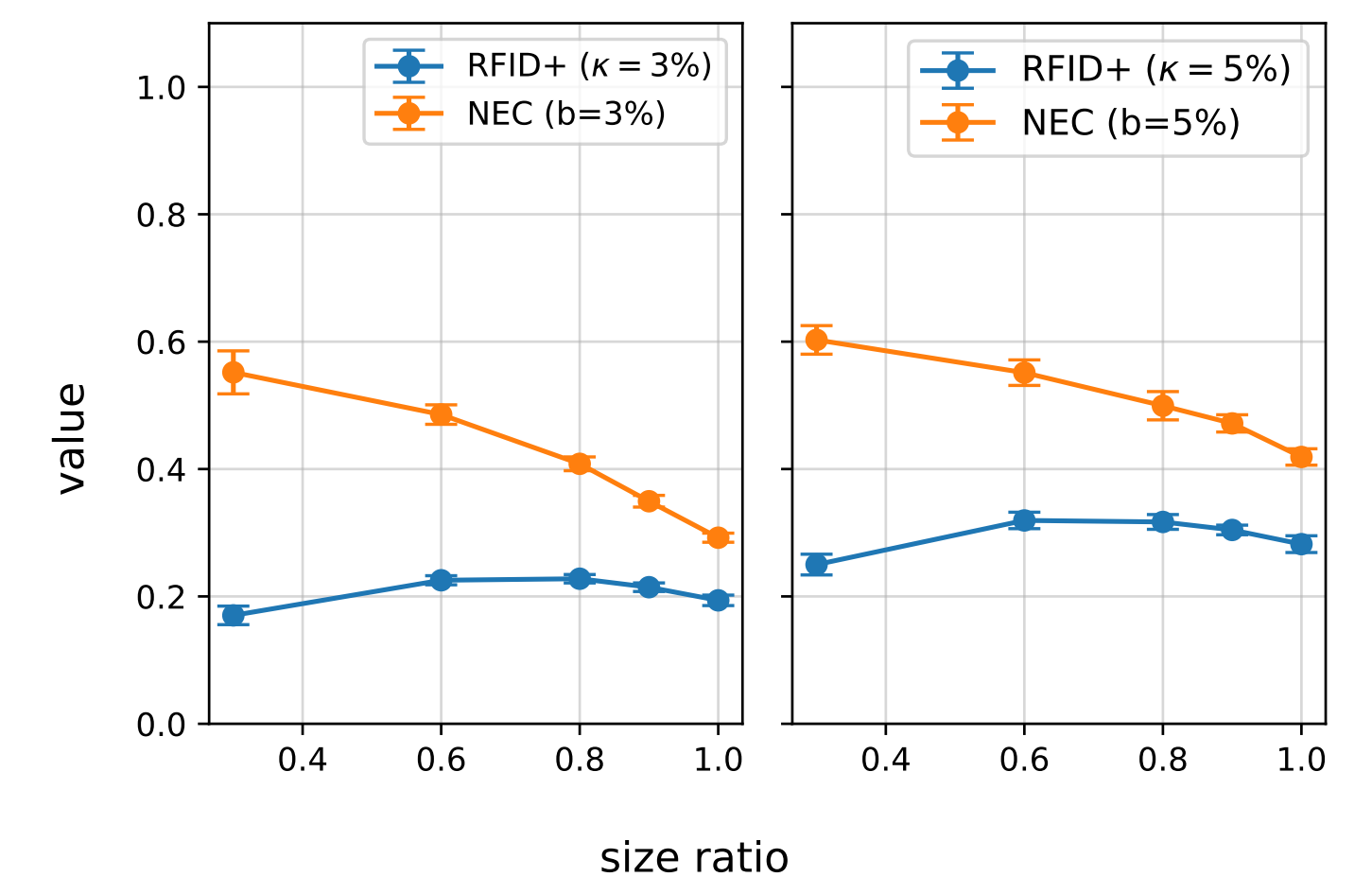


Figure 3. Our proposed Nec is sensitive to the number of irrelevant items in the explanation, whereas $RFid+$ is not.

RQ2: How good are SE-GNNs? [1]

We identified some architectural design choices favoring **un-faithfulness** and fixed them:

- **Hard Scores (HS):** give exact zero importance to information outside of R ;
- **Explanation Readout (ER):** aggregate only over R for the final prediction.

Table 3. Test set accuracy and faithfulness of some augmented SE-GNNs.

Dataset	BaMS		Motif2		Motif-Size		BBBBP	
	Acc	Faith	Acc	Faith	Acc	Faith	Acc	Faith
GSAT	100 ± 00	35 ± 03	92 ± 01	61 ± 01	90 ± 01	60 ± 02	79 ± 04	27 ± 08
GSAT + ER	100 ± 00	35 ± 03	92 ± 01	63 ± 01	90 ± 01	65 ± 01	80 ± 02	33 ± 04
GSAT + HS	98 ± 01	21 ± 06	53 ± 02	24 ± 05	54 ± 03	22 ± 05	71 ± 01	31 ± 09
GSAT + ER + HS	99 ± 01	24 ± 04	57 ± 04	37 ± 03	56 ± 07	29 ± 09	73 ± 02	33 ± 02
GISST	100 ± 00	25 ± 03	92 ± 01	53 ± 02	92 ± 00	50 ± 02	84 ± 03	23 ± 11
GISST + ER	-	-	-	-	-	-	85 ± 06	27 ± 06
GISST + HS	-	-	-	-	-	-	83 ± 05	19 ± 07
GISST + ER + HS	-	-	-	-	-	-	81 ± 07	15 ± 09
RAGE	96 ± 01	33 ± 05	83 ± 02	64 ± 04	74 ± 09	63 ± 07	82 ± 01	33 ± 04
RAGE + ER	96 ± 02	33 ± 02	85 ± 06	66 ± 03	71 ± 09	55 ± 07	84 ± 01	33 ± 05
RAGE + HS	97 ± 01	46 ± 03	85 ± 01	65 ± 02	78 ± 07	65 ± 09	84 ± 02	46 ± 02
RAGE + ER + HS	96 ± 01	46 ± 04	83 ± 04	64 ± 04	75 ± 08	62 ± 12	82 ± 01	43 ± 03

RQ3: Beyond subgraph-based explanations [2]

Theorem: Given a classifier g expressible as a purely existentially quantified first-order logic formula and a positive instance G of any size, then a Trivial Explanation for $g(G)$ is also a Prime Implicant explanation for $g(G)$.

- ⓘ Subgraph-based explanations are *optimal* for motif-based tasks;
- 😞 But we do not know when we are explaining motif-based tasks;
- 😞 Enhance standard SEGNNs with an interpretable side channel and let the optimization pick the best alternative (Occam's razor).

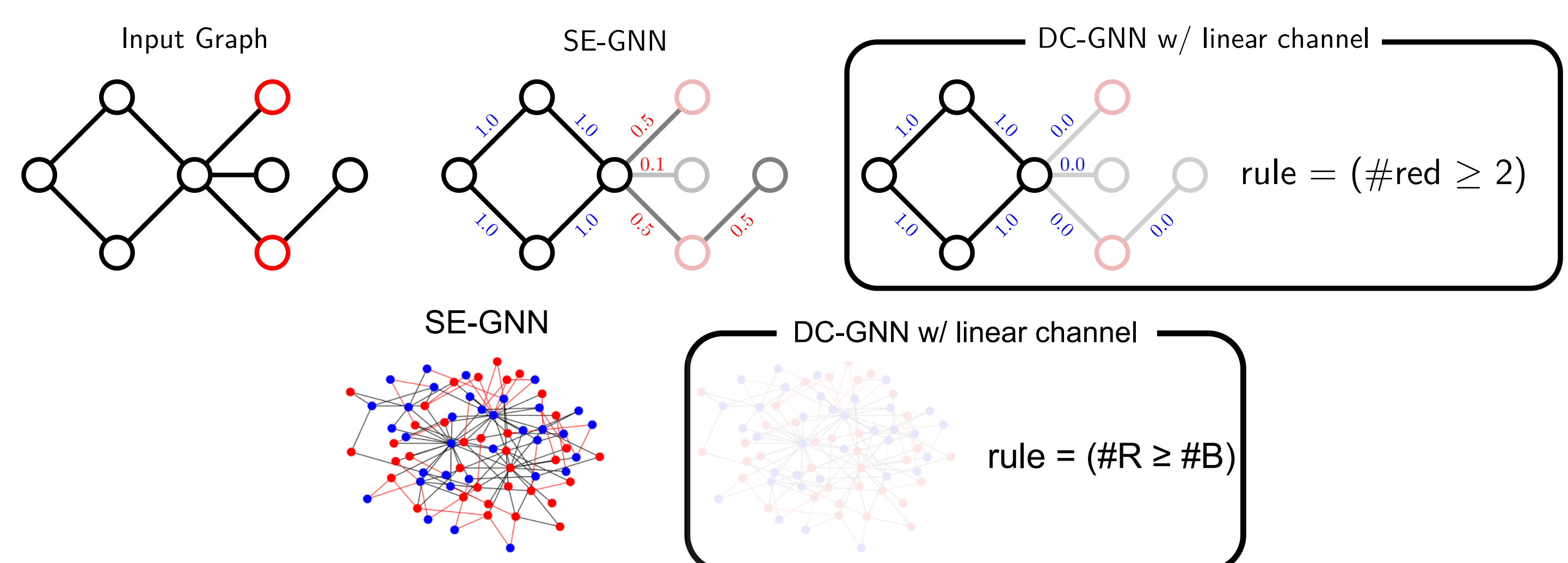


Figure 4. Examples of the proposed Dual-Channel SEGNN.

References

- [1] Steve Azzolin, Antonio Longa, Stefano Teso, and Andrea Passerini. Reconsidering faithfulness in regular, self-explainable and domain invariant GNNs. 2025.
- [2] Steve Azzolin, Sagar Malhotra, Andrea Passerini, and Stefano Teso. Beyond topological self-explainable gnn: A formal explainability perspective. 2025.
- [3] Marc Christiansen, Lea Villadsen, Zhiqiang Zhong, Stefano Teso, and Davide Mottin. How faithful are self-explainable gnn? 2023.
- [4] Zhong Li, Simon Geisler, Yuhang Wang, Stephan Günnemann, and Matthijs van Leeuwen. Explainable graph neural networks under fire. 2024.
- [5] Antonio Longa, Steve Azzolin, Gabriele Santin, Giulia Cencetti, Pietro Lio, Bruno Lepri, and Andrea Passerini. Explaining the explainers in graph neural networks: a comparative study. 2024.