# Reconsidering Faithfulness in Regular, Self-Explainable, and Domain Invariant GNNs

ICLR2025

Steve Azzolin, Antonio Longa, Stefano Teso, Andrea Passerini
University of Trento

📫 steve.azzolin@unitn.it

🏋 Graph Neural Networks (GNNs) rock at learning on graphs

😢 GNNs are black-box

🏋 Graph Neural Networks (GNNs) rock at learning on graphs

😢 GNNs are black-box

Post-hoc tools explain any given GNN (Amara et al., 2022; Kakkad et al., 2023; Longa et al., 2024)

Ante-hoc models provide an explanation <u>together with</u> the prediction (Miao et al., 2022; Kakkad et al., 2023; Longa et al., 2024)

How accurately does the explanation reflect the *reasoning* of the model?

How accurately does the explanation reflect the *reasoning* of the model?

👉 How **Faithful** is the explanation?

How accurately does the explanation reflect the *reasoning* of the model?

👉 How **Faithful** is the explanation?

Several different faithfulness metrics exist

🤔 which one to choose?

😨 what is faithfulness?

## Contributions

We propose to **reconsider faithfulness** from the following angles:

- ⚙️ How to compute faithfulness
- 🎯 How to promote faithfulness
- 📈 How does faithfulness affect model generalization

# How to compute faithfulness

We analyzed <u>seven</u> previous faithfulness metric and found that:

😨 Metric are not interchangeable[1] (Prop. 1)

😨 Some metrics do not encode the desired semantics (Prop. 2)

---
[1]Different metric yield different results

We analyzed <u>seven</u> previous faithfulness metric and found that:

😨 Metric are not interchangeable[1] (Prop. 1)

😨 Some metrics do not encode the desired semantics (Prop. 2)

😌 We propose a new necessity metric that penalizes overly large explanations (Prop. 3)

---

[1] Different metric yield different results

# How to promote faithfulness

🧐 Explanations of ante-hoc GNNs can be unfaithful (Christiansen et al., 2023)

We identified some architectural desiderata for faithfulness which ante-hoc GNNs do not implement:

🧐 Explanations of ante-hoc GNNs can be unfaithful (Christiansen et al., 2023)

We identified some architectural desiderata for faithfulness which ante-hoc GNNs do not implement:

😎 **Content Features**: the classifier should use raw input features, as opposed to the explanation extractor's embeddings

😎 **Explanation Readout**: the final global readout should run only over the explanation, as opposed to the entire graph

# How does faithfulness affect model generalization

Recent work in Invariant Learning for GNNs first **identifies an invariant subgraph**, and then **make predictions based on this subgraph only** (Chen et al., 2022; Gui et al., 2023)

The invariant subgraph plays the role of an explanation

😟 We show that extracting a domain invariant subgraph is not enough for a GNN to be truly domain invariant (Prop. 5)

😟 We show that extracting a domain invariant subgraph is not enough for a GNN to be truly domain invariant (Prop. 5)

> Unless the subgraph is also faithful, the information from the domain-dependent subgraph can still influence the prediction, thus preventing domain invariance

We proved that a GNN that fits well the ID data will fit well the OOD data if:

0  the explanation is in fact domain invariant

1  the explanation is faithful

# Reconsidering Faithfulness in Regular, Self-Explainable, and Domain Invariant GNNs

Steve Azzolin [1]   Antonio Longa [1]   Stefano Teso [1]   Andrea Passerini [1]

[1]University of Trento

UNIVERSITÀ DI TRENTO

## Motivation

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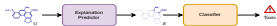


Figure 1. Pipeline of **Self-Explainable GNNs (SEGNNs)**.

⚠ How accurately does the explanation reflect the reasoning of the model?



👉 Measure the **faithfulness** of explanations, but multiple metrics exists
🙂 Which one to choose? 🤨 What is faithfulness?

## Our contribution 🎯

We propose to **reconsider faithfulness** from the following three angles

- How to compute faithfulness?
- How to promote faithfulness?
- How does faithfulness affect model generalization?

## How to compute faithfulness?

We abstract prior metrics into:

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

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

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

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

Table 1. SUF and NEC recover existing faithfulness metrics for appropriate choices of divergence $d$ and interventional distributions $p_d$ and $p_c$.

| Metric Estimates | | Divergence $d$ | Allowed changes |
|---|---|---|---|
| Uuf | Suf | $KL(p_d(\cdot \mid G) \parallel G')$ | zero out all irrelevant features |
| Fid | | $|p_d(\hat{y} \mid G) - p_d(\hat{y} \mid G')|$ | zero out all irrelevant features, delete all irrelevant edges |
| RFid | | | delete a random subset of irrelevant edges |
| FS | | $1\{p_d(\hat{y} \mid G) = p_d(\hat{y} \mid G')\}$ | multiply all relevant elements by relevance scores |
| Fid- | Nec | $|p_d(\hat{y} \mid G) - p_d(\hat{y} \mid G')|$ | zero out all relevant features, delete all relevant edges |
| RFid- | | | delete a random subset of relevant edges |
| PN | | $1\{p_c(\hat{y} \mid G) \neq p_d(\hat{y} \mid G')\}$ | multiply all relevant elements by relevance scores |

## How to compute faithfulness? (cont.)

Table 2. Model ranking and SUF results according to different $p_d$.

| Split | Model | Motif2 | |
|---|---|---|---|
| | | $p_d^{tr}$ | $p_d^{ID}$ |
| (I) | LECI | 1 (61 ±0.0) | 2 (62 ±0.0) |
| | GSAT | 2 (78 ±0.0) | 1 (84 ±0.0) |
| | CIGA | 3 (65 ±0.1) | 3 (73 ±0.4) |

🚫 **Metrics are not interchangeable**



Figure 2. Our proposed **Nec** is sensitive to the number of irrelevant items in the explanation, whereas **RFid-** is not.

🚫 Previous Necessity metrics **do not penalize useless explanations**

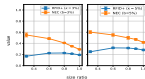🎯 We propose a **new necessity metric** penalizing overly large explanations

## How to promote faithfulness?

We identified some architectural design choices favoring **un-faithfulness** and fixed them:
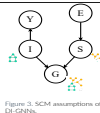
- **Hard Scores** (HS): give exact zero importance to information outside of $R$;
- **Explanation Readout** (ER): aggregate only over $R$ for the final prediction;
- **Content Features** (CF): feed the classifier raw features, not embeddings
- **Local Aggregations** (LA): non-local aggregations can create unwanted dependencies

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

| Dataset | BaMS | | Motif2 | | Motif-Size | | BBBP | |
|---|---|---|---|---|---|---|---|---|
| | Acc | Faith | Acc | Faith | Acc | Faith | Acc | Faith |
| GSAT | 100 ±0.0 | 35 ±1.2 | 92 ±0.1 | 61 ±1.9 | 90 ±0.0 | 60 ±0.0 | 79 ±0.7 | 27 ±1.9 |
| GSAT + ER | 100 ±0.0 | 35 ±2.7 | 92 ±1.2 | 63 ±2.3 | 90 ±0.0 | 65 ±0.0 | 80 ±0.8 | 33 ±3.3 |
| GSAT + HS | 98 ±1.7 | 21 ±3.3 | 53 ±3.2 | 26 ±3.3 | 54 ±0.0 | 22 ±2.1 | 71 ±0.6 | 25 ±5.2 |
| GSAT + ER + HS | 99 ±0.7 | 24 ±3.2 | 57 ±2.2 | 37 ±5.3 | 56 ±3.7 | 29 ±0.0 | 73 ±0.9 | 33 ±3.1 |
| GISST | 100 ±0.0 | 25 ±3.2 | 92 ±0.1 | 53 ±2.4 | 92 ±0.0 | 50 ±0.0 | 84 ±0.6 | 23 ±3.3 |
| GISST + ER | | | | | | | 85 ±0.8 | 27 ±2.9 |
| GISST + HS | | | | | | | 83 ±0.6 | 19 ±3.2 |
| GISST + ER + HS | | | | | | | 83 ±0.7 | 15 ±3.7 |
| RAGE | 96 ±0.0 | 33 ±2.2 | 83 ±0.8 | 64 ±2.8 | 74 ±0.0 | 63 ±0.2 | 75 ±1.0 | 28 ±3.4 |
| RAGE + ER | 96 ±0.0 | 33 ±3.3 | 83 ±0.6 | 66 ±0.1 | 71 ±3.5 | 54 ±0.0 | 84 ±2.4 | 33 ±3.6 |
| RAGE + HS | 97 ±1.1 | 46 ±2.6 | 85 ±0.5 | 72 ±0.3 | 70 ±2.5 | 65 ±0.2 | 84 ±1.7 | 46 ±3.4 |
| RAGE + ER + HS | 96 ±1.0 | 46 ±3.2 | 83 ±0.6 | 70 ±0.3 | 71 ±0.9 | 62 ±0.0 | 82 ±2.2 | 43 ±3.3 |

## How does faithfulness affect model generalization?

Recent work in Invariant Learning for GNNs first **identifies an invariant subgraph**, and then **make predictions based on this subgraph only** [1, 2].



Figure 3. SCM assumptions of DI-GNNs.

⚠ The invariant subgraph plays the role of an explanation

🙂 Unless the subgraph is also **faithful**, the information from the domain-dependent subgraph can still influence the prediction, thus **preventing domain generalization**

**Theorem 1.** Let $p_d$ be a deterministic DI-GNN with detector $f$ and classifier $g$, and $p^{\mathrm{ID}}(G, Y)$ and $p^{\mathrm{ood}}(G, Y)$ be the ID and OOD empirical distributions, respectively. Then:

$$\left| \mathbb{E}_{(G,y) \sim p^{\mathrm{ID}}}[p_d(y \mid G)] - \mathbb{E}_{(G,y) \sim p^{\mathrm{ood}}}[p_d(y \mid G)] \right| \quad (1)$$
$$\leq \mathbb{E} \left[ k_1(\lambda_{\mathrm{corr}}^{\mathrm{id}} + \lambda_{\mathrm{corr}}^{\mathrm{ood}}) + k_2(\lambda_{\mathrm{suf}}^{\mathrm{id}} + \lambda_{\mathrm{suf}}^{\mathrm{ood}}) + (\lambda_{\mathrm{suf}}^{\mathrm{id}} + \lambda_{\mathrm{suf}}^{\mathrm{ood}}) \right]$$

🎯 A DI-GNN that fits the ID data well will fit the OOD data well if:

- $R$ is domain-invariant (low $\lambda_{\mathrm{corr}}$ and $\lambda_{\mathrm{suf}}$)
- highly sufficient (low $\lambda_{\mathrm{suf}}$)

🎯 **Correlate** the difference in average prediction's **likelihood** between ID and OOD data, and the sum of the **degree of domain-invariance** and **faithfulness**
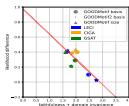


Figure 4. Likelihood, faithfulness, and domain-invariance are correlated.

## References

[1] Yongqiang Chen, Yonggang Zhang, Yatao Bian, Han Yang, MA Kaili, Binghui Xie, Tongliang Liu, Bo Han, and James Cheng.
Learning causally invariant representations for out-of-distribution generalization on graphs.
Advances in Neural Information Processing Systems, 35:22131–22148, 2022.

[2] Shurui Gui, Meng Liu, Xiner Li, Youzhi Luo, and Shuiwang Ji.
Joint learning of label and environment causal independence for graph out-of-distribution generalization.
In Thirty-seventh Conference on Neural Information Processing Systems, 2023.

✉ steve.azzolin@unitn.it

# References

Kenza Amara, Zhitao Ying, Zitao Zhang, Zhichao Han, Yang Zhao, Yinan Shan, Ulrik Brandes, Sebastian Schemm, and Ce Zhang. 2022. GraphFramEx: Towards Systematic Evaluation of Explainability Methods for Graph Neural Networks. In *Learning on Graphs Conference*. PMLR, 44–1.

Yongqiang Chen, Yonggang Zhang, Yatao Bian, Han Yang, MA Kaili, Binghui Xie, Tongliang Liu, Bo Han, and James Cheng. 2022. Learning causally invariant representations for out-of-distribution generalization on graphs. *Advances in Neural Information Processing Systems* 35 (2022), 22131–22148.

Marc Christiansen, Lea Villadsen, Zhiqiang Zhong, Stefano Teso, and Davide Mottin. 2023. How Faithful are Self-Explainable GNNs? *arXiv preprint arXiv:2308.15096* (2023).

Shurui Gui, Meng Liu, Xiner Li, Youzhi Luo, and Shuiwang Ji. 2023. Joint Learning of Label and Environment Causal Independence for Graph Out-of-Distribution Generalization. *arXiv preprint arXiv:2306.01103* (2023).

Jaykumar Kakkad, Jaspal Jannu, Kartik Sharma, Charu Aggarwal, and Sourav Medya. 2023. A Survey on Explainability of Graph Neural Networks. *arXiv preprint arXiv:2306.01958* (2023).

Antonio Longa, Steve Azzolin, Gabriele Santin, Giulia Cencetti, Pietro Lio, Bruno Lepri, and Andrea Passerini. 2024. Explaining the Explainers in Graph Neural Networks: a Comparative Study. *ACM Comput. Surv.* (2024). https://doi.org/10.1145/3696444

Siqi Miao, Mia Liu, and Pan Li. 2022. Interpretable and generalizable graph learning via stochastic attention mechanism. In *International Conference on Machine Learning*. PMLR, 15524–15543.