


Faithful and Accurate Self-Attention Attribution for Message Passing Neural Networks via the Computation Tree Viewpoint

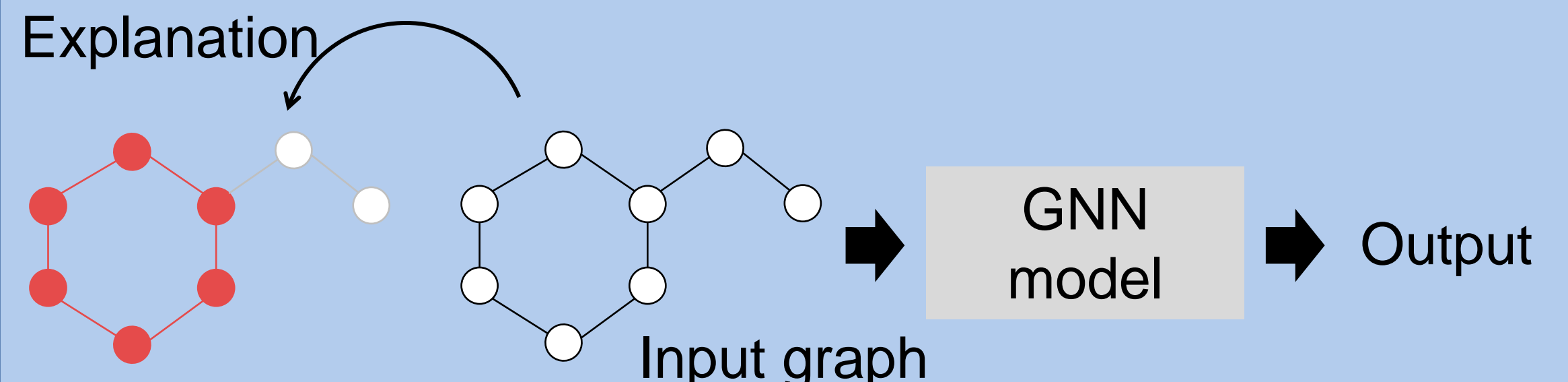
1. Vast discussion on Attention in CV & NLP

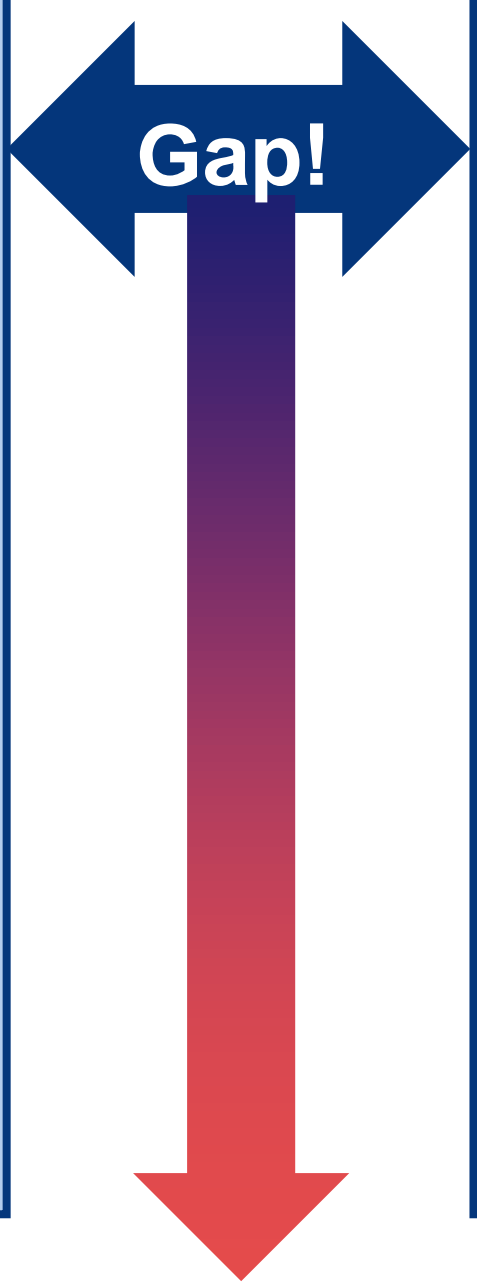
- Attention is already used as explanations of transformer models in computer vision: Rollout [1], Chefer et al. [2 & 3]
- A long discussion of attention-based interpretation in natural language processing (see [4] for full survey)



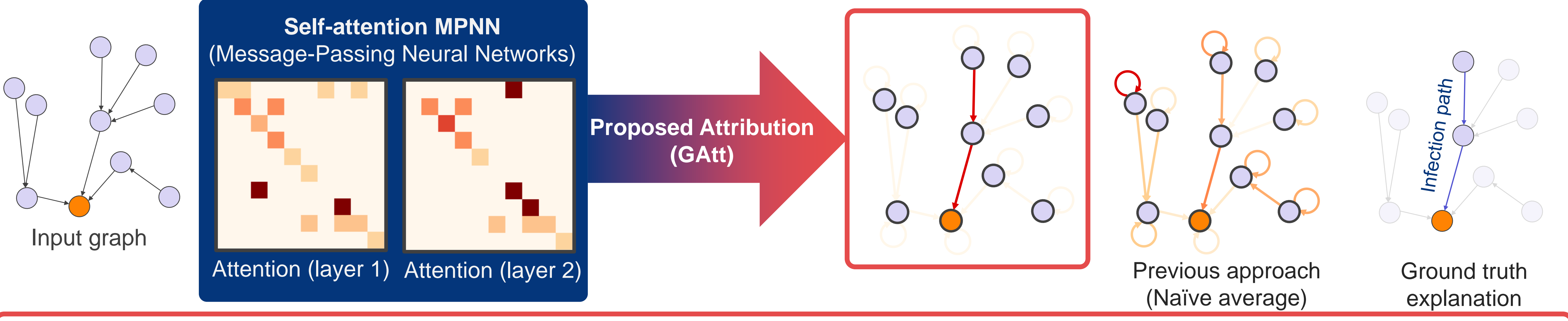
2. Underdeveloped topic in explaining GNNs

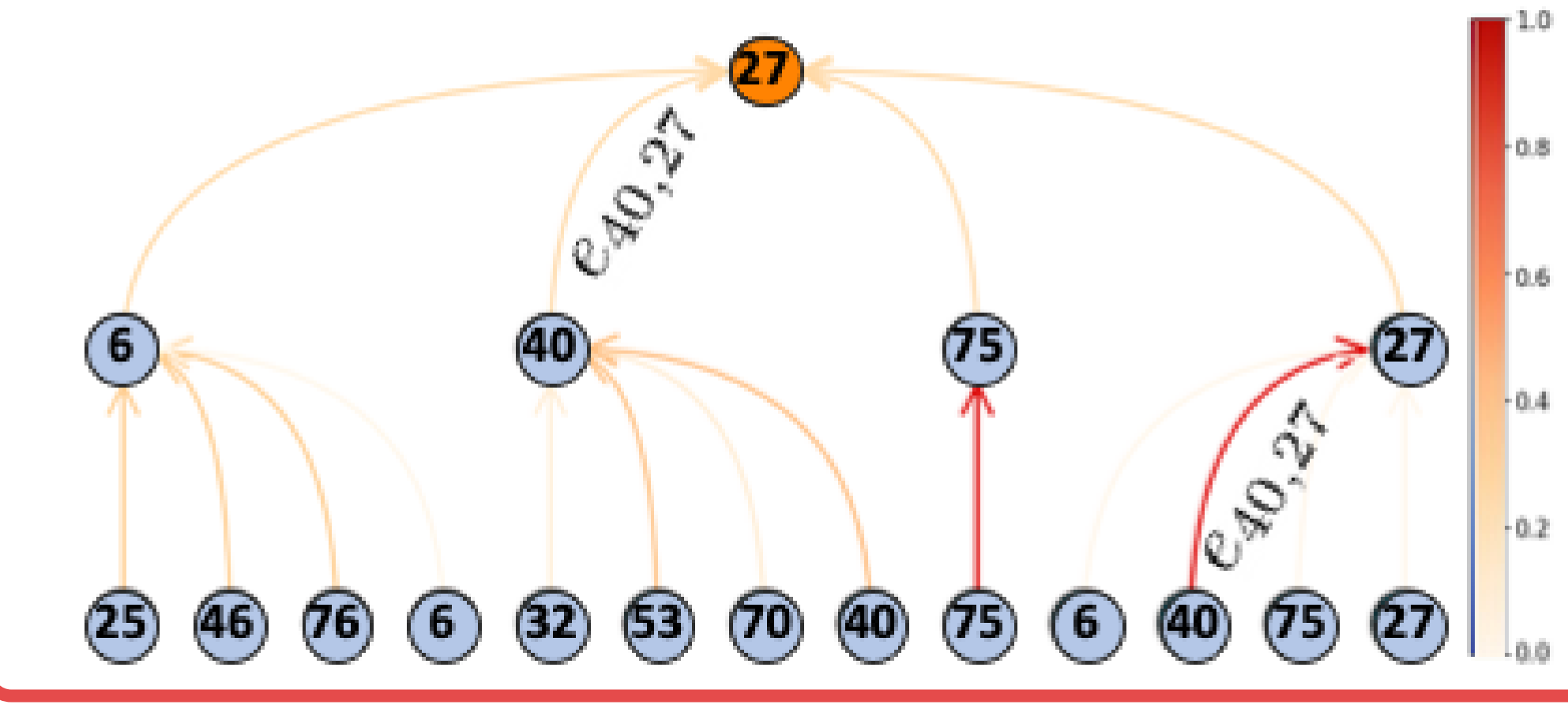
- XAI in GNNs is actively studied in recent years. Explanation methods usually aim to find the most relevant part of the input (attribution)
- However, there has been nearly no discussion on the potential of attention as explanation in GNNs





Our work is one of the first to seriously discuss the **potential of attention weights** in **MPNNs** as explanations.





Design Philosophy of GAtt:

Computation Tree Viewpoint of Edge Attribution Calculation from Attention

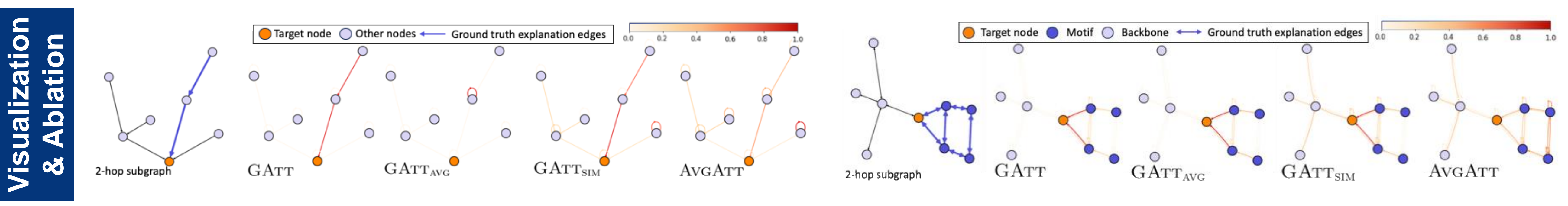
- Proximity effect:** Edges within closer proximity to the target node tend to highly impact the model's prediction compared with distant edges, since they are likely to appear more frequently in the computation tree. **Need to sum all occurrences of an edge!**
- Contribution adjustment:** The contribution of an edge in the computation tree should be **adjusted by its position (i.e., other edges in the path towards the root)**.

Faithfulness evaluation	Dataset		2-layer GAT/GATv2			3-layer GAT/GATv2		
			GATT	AVGATT	Random	GATT	AVGATT	Random
Cora	Δ_{PC}		0.8468/0.1040	0.1764/0.0121	-0.0056/-0.0036	0.8642/0.1696	0.0967/0.0168	0.0045/0.0045
	Δ_{NE}		0.7112/0.0930	0.1526/0.0100	-0.0076/0.0019	0.7690/0.1664	0.0859/0.0186	0.0040/0.0037
	Δ_P		0.9755/0.9623	0.7251/0.6226	0.4389/0.4891	0.9875/0.9966	0.7075/0.8897	0.5235/0.6107

- Faithfulness:** How much does the attribution actually reflect the behavior of the underlying model?
- Results (7 datasets, 3 measurements, 2/3 layer **GAT** [5] / **GATv2** [6] / **SuperGAT** [7]) show clear superiority.

Accuracy evaluation	Model	Dataset	GATT	AVGATT	SA	GB	IG	GNNE _x	PGE _x	GM	FDnX	Random
GAT		BA-Shapes	0.9591	0.7977	0.9563	0.6231	0.6231	0.8916	0.8289	0.5316	0.9917	0.4975
		Infection	0.9976	0.8786	0.8237	0.8949	<u>0.9472</u>	0.9272	0.7173	0.6859	0.6574	0.4811

- Accuracy:** How much does the attribution reveal the ground truth explanation?
- Results (2 datasets, 1 attention baseline, **7 post-hoc explanation baselines**) show clear competent performance.



GAtt shows superior empirical performance on a wide variety of datasets, measures, and attention-based MPNN models.