# Rethinking Attention-Model Explainability through Faithfulness Violation Test

Yibing Liu [1]   Haoliang Li [1]   Yangyang Guo [2]   Chenqi Kong [1]   Jing Li [3]   Shiqi Wang [1]

## Abstract

Attention mechanisms are dominating the explainability of deep models. They produce probability distributions over the input, which are widely deemed as feature-importance indicators. However, in this paper, we find one critical limitation in attention explanations: weakness in identifying the polarity of feature impact. This would be somewhat misleading – features with higher attention weights may not faithfully contribute to model predictions; instead, they can impose suppression effects. With this finding, we reflect on the explainability of current attention-based techniques, such as Attention ⊙ Gradient and LRP-based attention explanations. We first propose an actionable diagnostic methodology (henceforth *faithfulness violation test*) to measure the consistency between explanation weights and the impact polarity. Through the extensive experiments, we then show that most tested explanation methods are unexpectedly hindered by the faithfulness violation issue, especially the raw attention. Empirical analyses on the factors affecting violation issues further provide useful observations for adopting explanation methods in attention models.

## 1. Introduction

Recent advances in attention models (Bahdanau et al., 2015; Vaswani et al., 2017) have raised a pressing concern for the explainability of their black-box nature. Prior efforts target at unraveling puzzles in inherent attention explanations, such as whether attention is directly explainable (Jain & Wallace, 2019; Serrano & Smith, 2019; Wiegreffe & Pinter, 2019; Brunner et al., 2020), or how attention behaves (Clark et al., 2019; Bai et al., 2021; Michel et al., 2019; Mohankumar et al., 2020). In contrast, another line of work strives to

[1]City University of Hong Kong, Hong Kong [2]National University of Singapore, Singapore [3]The Hong Kong Polytechnic University, Hong Kong. Correspondence to: Haoliang Li <haoliang.li1991@gmail.com>.
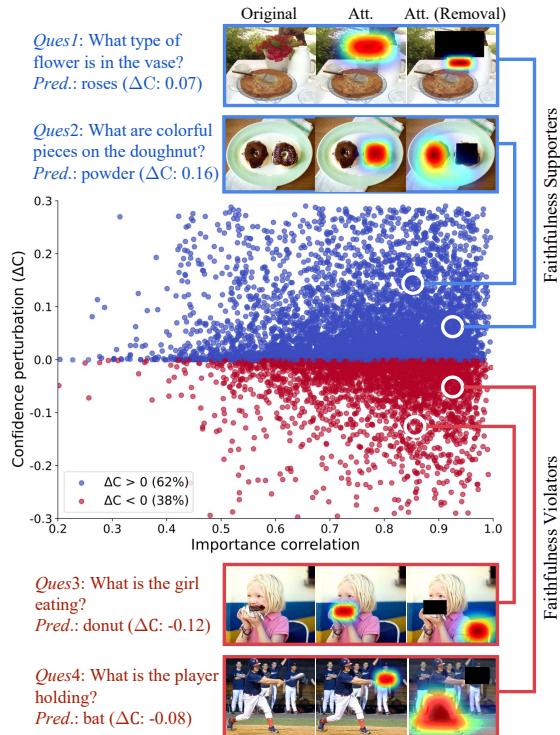
*Figure 1.* Data map for attention faithfulness on BUTD model over VQA v2. The $y$-axis shows the confidence perturbation $\Delta C$ on the original predicted label (*e.g.*, *donut*) after removing top-10% image regions *w.r.t* attention weights. $\Delta C < 0$ indicates confidence increases. The *Red* points represent the ones violating *polarity consistency* since image regions with the largest positive weights unexpectedly suppress model predictions, *i.e.*, $\Delta C < 0$. The $x$-axis shows the *importance correlation* for attention weights.

devise better explanation methods for attention models (Sundararajan et al., 2017; Smilkov et al., 2017; Shrikumar et al., 2017; Chefer et al., 2021a; Hao et al., 2021; Kobayashi et al., 2020). However, in this work, we find that they all lack proper evaluation in *explanation faithfulness* (Jacovi & Goldberg, 2020) – how accurately an explanation reflects the reasons behind a model prediction.

We argue that a faithful explanation should at least possess two properties: 1) *importance correlation*: the magnitude of explanation weights should precisely reflect the importance of input components; 2) *polarity consistency*: the sign of

explanation weights should correctly indicate the polarity of input impact, *i.e.*, contribution or suppression effects to the model predictions. Existing faithfulness evaluations, however, merely focus on the first property without examining the polarity consistency (Jain & Wallace, 2019; Wiegreffe & Pinter, 2019; Serrano & Smith, 2019; Atanasova et al., 2020; Moradi et al., 2021; Chrysostomou & Aletras, 2021). While somewhat dissatisfying, this can be heavily misleading: in cases where an explanation confuses the polarity but well correlates with importance, we may mistakenly identify inputs (*e.g.*, image regions) with large suppression effects as the reason for model predictions.

To illustrate this point, we build a dataset map (Swayamdipta et al., 2020) in Figure 1 for attention faithfulness evaluation over the Visual Question Answering (VQA) task, which distinguishes between the samples that violate polarity consistency (*Red* points) and the ones support (*Blue* points). As can be seen, in numerous appealing cases where even attention weights precisely reflect importance of image regions (*i.e.*, high importance correlation), highlighted regions may not faithfully contribute to the predictions; instead, they impose suppression effects, *i.e.*, violating polarity consistency. This hence would be fallacious if we believe such highlighted image regions help the model prediction. As such, in this paper, we argue that polarity consistency is essential to be a criterion for faithful explanations, and a re-assessment for explanation faithfulness is imperative.

To this end, we propose an actionable diagnostic methodology, *violation test*, to evaluate the under-explored property, polarity consistency, in faithfulness. In particular, this violation test would take the largest explanation weight (in absolute value) as the representative for each explanation sample, and then examine the correctness of its corresponding polarity, *i.e.*, whether denoting the actual direction of input impact. In response, if the polarity is correct, such an explanation sample would be marked as the supporter for faithfulness; otherwise, the violator. We should expect an explanation method passes the test with as few violations as possible, if it faithfully reflects model decisions. In addition, we also compare our proposed violation test with existing faithfulness evaluation metrics (Section 5.1). Through the results, we show that they all lack the ability in examining polarity consistency, which further highlights the necessity of our violation test for faithfulness evaluation. Note that our test method is also applicable to any explanation method without modifications.

With the realization of the weakness in existing faithfulness evaluation, we instantiate our analysis on nine widely utilized explanation methods, *e.g.*, Attention Rollout (Abnar & Zuidema, 2020). In our broad experimental sweep over six tasks and nine datasets, we find that most analyzed explanation methods are unexpectedly hindered by the faithfulness

violation issue (Section 5.2). Specifically, the single-polarity methods (*e.g.*, Raw Attention) perform the worst, while the form of gradient-based attention explanations can largely alleviate violation issue. With these findings, we further study the factors leading to the violation issue (Section 5.3). Based on our experiments, we show that the complexity in model architectures and the capability to identify polarity can be the main factors.

To sum up, this work makes the following contributions:

- We study the under-explored dimension, *polarity consistency*, in the faithfulness evaluation. We propose a violation test to re-assess the faithfulness of explanation methods. By comparison with existing faithfulness evaluation metrics, we show the indispensability of our proposed violation test.

- We conduct extensive experiments with three groups of explanation methods across tasks, datasets, and model architectures. We find, consistently, most methods tested are limited by the faithfulness violation issue. In particular, raw attention explanations perform the worst, while the form of attention ⊙ gradient can greatly alleviate the violation issue and even pass the violation test in some special cases.

- We investigate the dominant factors leading to the violation issue. The results suggest that the complexity in model architectures and the capability of identifying polarity play important roles. All code for model implementation and data processing is made avaliable[1].

## 2. Related Work

### 2.1. Model Explainability

Explaining the internal behavior of black-box neural networks has been extensively studied in recent years. They can be grouped into several families, including but not limited to gradient-based methods (Shrikumar et al., 2017; Sundararajan et al., 2017; Binder et al., 2016; Smilkov et al., 2017), perturbation-based explanations (Zeiler & Fergus, 2014; Zintgraf et al., 2017; Fong et al., 2019), and Shapley values (Lundberg & Lee, 2017; Chen et al., 2019).

Specific to attention models, a canonical explanation tool is attention mechanisms (Bahdanau et al., 2015; Vaswani et al., 2017). Since attention calculates a probability distribution over the input, attention weights are often compared to human vision (Rensink, 2000) and regarded as a kind of importance indicator (Mnih et al., 2014; Luong et al., 2015; Yang et al., 2016). In line with this, many works hence develop explanation methods by combining attention weights with additional attributions, such as the gradients (Chefer

---

[1]https://github.com/BierOne/Attention-Faithfulness

et al., 2021a; Hao et al., 2021), layer-wise relevance (Voita et al., 2019; Chefer et al., 2021b), or the norm of input vectors (Kobayashi et al., 2020). On the other end of the spectrum, many studies also seek to understand the behavior of attention. For example, Mohankumar et al. (2020) explored why attention fails to explain LSTM models and pointed out that the similar hidden states in the encoder restrict the significance of attention weights. In contrast, another group of work investigates what patterns attention often looks at (Clark et al., 2019; Bai et al., 2021). Unlike them, in this paper, we provide insights for attention-based explanations from the perspective of polarity consistency.

This work is also broadly connected to the evaluation of saliency methods (Adebayo et al., 2018; Tomsett et al., 2020; Srinivas & Fleuret, 2021). Contrary to these works mostly focusing on how explanations can be arbitrarily manipulated, our polarity consistency view centers on evaluating polarity indications of explanation weights, which sheds light on another research direction in XAI.

### 2.2. Faithfulness in Attention Models

An explanation satisfying faithfulness should at least possess two properties: 1) importance correlation: precisely reflecting the importance of each input feature, and 2) polarity consistency: correctly identifying the impact polarity of features. Despite recent efforts on faithfulness evaluation where many metrics have been proposed (DeYoung et al., 2020; Atanasova et al., 2020), they merely focus on the first property without examining the polarity consistency.

The first attempt to explore faithfulness in attention models is Jain & Wallace (2019); Serrano & Smith (2019), where they found that large attention weights do not necessarily indicate high importance on model predictions and thus drew the conclusion that attention is not faithful. An immediate response (Wiegreffe & Pinter, 2019), however, argues their experiments and presents a different claim that attention can provide faithful explanations in certain cases. Furthermore, another work compares attention weights with other explanations, such as LIME (Ribeiro et al., 2016; Thorne et al., 2019) and Shapley values (Ethayarajh & Jurafsky, 2021). In addition to the exploration of faithfulness, some studies also propose to improve faithfulness of attention-based explanations on the specific domains (Moradi et al., 2021; Chrysostomou & Aletras, 2021). However, their faithfulness measurements are all performed on the importance correlation of attention-based explanations, which heavily limit the reliability of the results.

### 3. Faithfulness Violation Test

The problem we are investigating is to what extent an explanation method can correctly identify the impact polarity of

*Table 1.* Task and dataset statistics. We provide more details for each dataset in Appendix A.

| Task | Dataset | # Classes | # Train |
|---|---|---|---|
| Sentiment Analysis | SST | 5 | 6,355 |
| | Yelp | 5 | 650,000 |
| Topic Classification | AGNews | 4 | 101,998 |
| | 20News | 20 | 10,833 |
| Paraphrase Detection | QQP | 2 | 327,460 |
| Natural Language Inference | SNLI | 3 | 549,367 |
| Question Answering | bAbI-1 | 20 | 10,000 |
| Visual Question Answering | VQA 2.0 | 3130 | 411,272 |
| | GQA | 1842 | 943,000 |

input features. Without loss of generality, suppose we have an input $x = \{x_i\}_i^N$, where each $x_i$ represents a feature or token. A model is represented as a function $f(\cdot)$, predicting the class $\hat{y}$ for the input: $\hat{y} = \mathrm{argmax}_y f(x)_y$. Let $w$ be an explanation method, which assigns weights for each token in the input: $e_i = w(x_i), i \in [1, N]$.

The intuition behind faithfulness violation test is that features with the largest explanation weights (in absolute values) should at least demonstrate the consistent polarity, *e.g.*, largest positive weights ought to indicate contribution effects instead of suppression. In this sense, our first step is to find the most influential feature $x^*$ (denoted by the explanation method) in the input: $x^* = \mathrm{argmax}_{x_i \in x} ||w(x_i)||$. Next, to measure the actual impact of $x^*$, we perform a well-established approach – perturbation examination (Li et al., 2016; Hooker et al., 2019; DeYoung et al., 2020) – replacing the certain input features ($x^*$) and calculating the confidence perturbation in model predictions:

$$\Delta \mathrm{C}(x, x^*) = f(x)_{\hat{y}} - f(x \backslash x^*)_{\hat{y}}. \quad (1)$$

In particular, if the explanation weights $w(x^*)$ are highly positive, the confidence degradation ($\Delta \mathrm{C}(x, x^*) > 0$) should be observed after the feature replacement. Should the confidence upgrades, however, we can infer that such an explanation violates the polarity consistency. As such, the faithfulness violation is measured via

$$\mathrm{Violation} = \mathbb{1}_{\mathrm{sign}(w(x^*) \cdot \Delta \mathrm{C}(x, x^*)) < 0} \quad (2)$$

### 4. Experimental Setup

#### 4.1. Tasks and Datasets

Our goal is to explore the faithfulness violation issue of attention-based explanations. Following prior studies (Jain & Wallace, 2019; Mohankumar et al., 2020), we conduct extensive experiments on six exemplar tasks, for which

*Table 2.* The most common explanation techniques for attention models are analyzed in this paper. $\alpha$ denotes the raw attention weights. $||v(x)||$ represents the norm of transformed input vectors in attention module. $R^\alpha$ is the Layer-wise Relevance Propagation (LRP) scores for the corresponding attention head. Appendix B provides an in-depth description for these techniques.

| Method | Denoted | Basis |
|---|---|---|
| *Generic attention-based explanation methods* | | |
| Inherent Attention Explanation | RawAtt | $\alpha$ |
| Attention $\odot$ Gradient | AttGrad | $\alpha \odot \nabla\alpha$ |
| Attention $\odot$ InputNorm | AttIN | $\alpha \odot ||v(x)||$ |
| *Transformer-based explanation methods* | | |
| Partial LRP | PLRP | $R^\alpha$ |
| Attention Rollout | Rollout | $\alpha$ |
| Transformer Attention Attribution | TransAtt | $\nabla\alpha \odot R^\alpha$ |
| Generic Attention Attribution | GenAtt | $\alpha \odot \nabla\alpha$ |
| *Gradient-based attribution methods* | | |
| Input $\odot$ Gradient | InputGrad | $x \odot \nabla x$ |
| Integrated Gradients | IG | $x \odot \nabla x$ |

attention models are widely applied. Besides, to ensure the reliability of our experimental results, we select nine datasets across different training data scales, class numbers, and learning targets. We report the setup of them in Table 1.

## 4.2. Models

In line with previous studies (Serrano & Smith, 2019; Jain & Wallace, 2019; Atanasova et al., 2020), we consider two commonly employed base architectures, LSTM (Hochreiter & Schmidhuber, 1997) and CNN (LeCun et al., 2010), for the feature extraction on the standard NLP tasks, *e.g.*, *Sentiment Analysis*. Since our goal is to analyze attention-based explanations, we incorporate two kinds of attention mechanisms, additive attention (*TanhAtt*) (Bahdanau et al., 2015) and scaled dot-product attention (*DotAtt*) (Vaswani et al., 2017), for the above architectures as well. As such, we experiment with four general attention models here: **LSTM+TanhAtt**, **LSTM+DotAtt**, **CNN+TanhAtt**, and **CNN+DotAtt**, all of which is initialized with the GloVe (Pennington et al., 2014) embeddings.

We utilize two typical transformer-based models, **Visual-BERT** (Li et al., 2020) and **LXMERT** (Tan & Bansal, 2019), on the VQA task since they are representatives of self-attention and cross-modal attention mechanisms, respectively. In addition, we also use two widely adopted general attention models, **StrongBaseline** (Kazemi & Elqursh, 2017) and **BUTD** (Anderson et al., 2018), on the VQA task for a thorough comparison. Another difference between these four models is that StrongBaseline performs on the grid-based image features from ResNet-101 (He et al., 2016), while others employ the object-level image features from Faster R-CNN (Ren et al., 2015).

*Table 3.* The faithfulness evaluation metrics utilized in this paper. Appendix C provides more details for these metrics.

| Metric | Denoted | Measurement |
|---|---|---|
| AUC-TP | A $\downarrow$ | Performance perturbation |
| Sufficiency | S $\downarrow$ | Confidence perturbation |
| Comprehensiveness | C $\uparrow$ | Confidence perturbation |
| Rank Correlation | R $\uparrow$ | Overall confidence variance |
| Ours Violation Test | V $\downarrow$ | Confidence perturbation polarity |

## 4.3. Explanation Methods

We mainly study the attention model explainability from three groups of explanation methods – *generic attention-based*, *transformer-based*, and *gradient-based* attribution approaches. In particular, generic attention-based methods are suitable for all kinds of attention models, while transformer-based ones are exclusive for transformer architectures. We elaborate the utilized methods in Table 2.

Starting with generic attention-based explanation methods, we employ Inherent Attention Explanation (RawAtt) (Jain & Wallace, 2019), Attention $\odot$ Gradient (AttGrad) (Serrano & Smith, 2019; Chrysostomou & Aletras, 2021; Hao et al., 2021), and Attention $\odot$ InputNorm (AttIN) (Kobayashi et al., 2020) for all models. Note that in transformer-based architectures, these methods are implemented based on the last attention layer's output following (Chefer et al., 2021a).

Regarding the exclusive architecture of transformers (*e.g.*, multi-head attention), we particularly adopt four explanation methods: Partial Layer-wise Relevance Propagation (PLRP) (Voita et al., 2019), Attention Rollout (Abnar & Zuidema, 2020), Transformer Attention Attribution (TransAtt) (Chefer et al., 2021b), and Generic Attention Attribution (GenAtt) (Chefer et al., 2021a).

Additionally, in our experiments, we also select two representatives – Input $\odot$ Gradient (InputGrad) (Shrikumar et al., 2017; Ancona et al., 2018) and Integrated Gradients (IG) (Sundararajan et al., 2017) – from gradient-based attribution methods for a further comparison.

## 4.4. Faithfulness Evaluation Metrics

We compare our violation test with existing faithfulness evaluation metrics, while they mainly focus on the importance correlation of explanation methods. Concretely, we experiment with four widely adopted metrics as following.

**AUC on Threshold Performance (AUC-TP)** (Atanasova et al., 2020). This metric calculates the AUC score based on the performance perturbation *w.r.t* the feature replacement of top-0, 5, 10, 20, 30, ..., 90% tokens in the input.

**Sufficiency** (DeYoung et al., 2020) measures whether explanations can identify important features, which are adequate

for a model to remain confidence on its original predictions. It is calculated as the average of $\Delta C$ by sequentially retaining the features in top-5, 10, 20, 50% sparsity levels.

**Comprehensiveness** (DeYoung et al., 2020) evaluates if the features assigned with lower weights are unnecessary for the predictions. In analogous to Sufficiency, we compute it as the average of $\Delta C$ via sequentially removing the features in the top-5, 10, 20, 50% sparsity levels.

**Rank Correlation (Absolute)** (Jain & Wallace, 2019; Selvaraju et al., 2020) directly calculates the perturbation-based prediction variance for each token, and measures its relationship with the magnitude of explanation weights.

### 4.5. Replacement Functions

In the light of out-of-distribution (OOD) problems in perturbation examination (Hooker et al., 2019), we utilize three replacement functions as suggested by Hase et al. (2021). We provide a brief overview of the involved functions in the following. Note that in this work, all replacements are performed on the token level, *i.e.*, text token or image region.

**Slice Out** (DeYoung et al., 2020) directly removes specified tokens from the original input, thereby resulting in a shorter length of the new input sequence.

**Attention Mask** (Serrano & Smith, 2019) sets the attention weights for removed tokens to 0, such that the removed parts will not be forwarded in the network anymore.

**MASK Token** (Li et al., 2016; Sundararajan et al., 2017) simply replaces the words with the MASK token. If the token is an image region, we would set it to the zero vector.

In the following sections, all the results are averaged by using the above three replacement functions.

## 5. Experimental Results

Faithfulness violation indicates that an explanation does not truly reflect the impact of input components. In this work, we mainly study the faithfulness violation issue in terms of the polarity consistency, *i.e.*, whether an explanation can correctly identify the impact polarity of features. In particular, we consider the context of attention models due to their wide applications. As described in the previous section, we explore 9 explanation methods across 8 popular model architectures over 6 tasks. To sum up, our exploration is to answer the following research questions:

- **RQ1**: Why we need faithfulness violation test? Can existing evaluation metrics reflect the violation issue?

- **RQ2**: How existing methods perform on faithfulness? Can any method pass the faithfulness violation test?
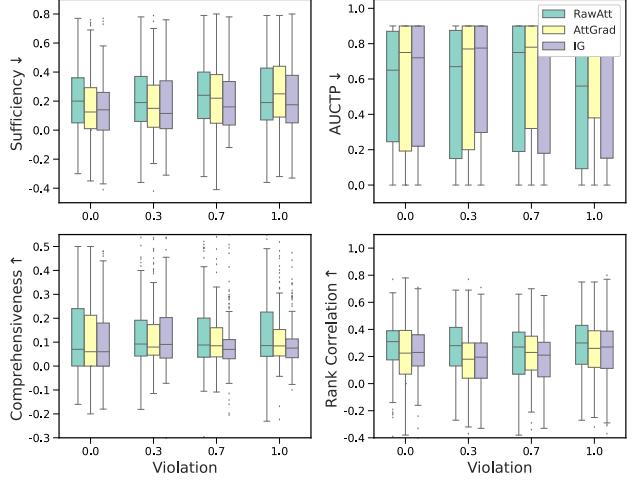


*Figure 2.* Faithfulness violation levels ($x$-axis) in correlation to the existing faithfulness evaluation metrics ($y$-axis). The violation level is obtained by the average of our violation results, which denotes worse faithfulness as it increases. The presented explanation methods (RawAtt, AttGrad, and IG) are based on the LXMERT model over the GQA dataset.

- **RQ3**: What factors affect the violation issue?

Regarding the first question, we begin by comparing our violation test with existing faithfulness evaluation metrics. We show that all of these metrics are incapable of detecting faithfulness violations. In addition, we further verify the function of our violation test, where the results demonstrate the capability of the test method in detecting violation behavior. In line with such experiments, for the second question, we conduct a sanity faithfulness evaluation for existing explanation methods. Unexpectedly, the results suggest that almost all approaches are hindered by the violation issue. In particular, the single-polarity methods (*e.g.*, RawAtt and AttIN) perform worst, while the gradient-based explanations (*e.g.*, GenAtt) could greatly alleviate the violation issue. Finally, for the third question, we study the factors leading to the above issues. Through the experiments, we find that the complexity of model architectures and the ability to identify polarity can be the dominant factors.

### 5.1. Why violation test is necessary? (RQ1)

**Comparison with existing evaluation metrics.** We first compare the existing evaluation metrics with our proposed violation test. Concretely, we randomly select the same number of data samples at different violation levels, and then calculate their corresponding performance distribution for each evaluation metric (see Figure 2). To our surprise, it can be found that at different violation levels, the performance of explanation methods on the same evaluation metric is
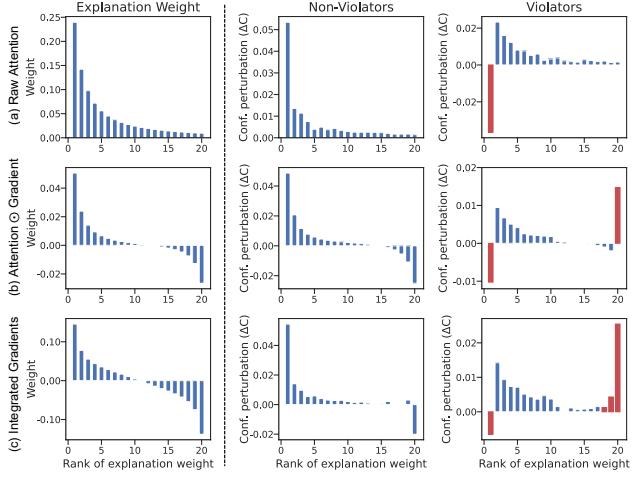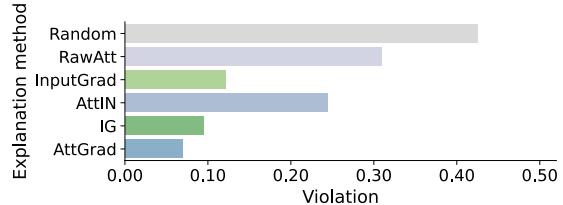
*Figure 3.* Comparison between the identified faithfulness non-violators (middle figures) and violators (right ones). The leftmost figures exhibit ordered explanation weights at each rank. The rest figures represent the correlation between explanation weights and perturbation-based input impact ($\Delta C$). *Red* bars denote the faithfulness violation in polarity consistency, *i.e.*, the polarity of explanation weight is opposite with the $\Delta C$. The results are presented on the LXMERT model over the GQA dataset.

*Figure 4.* Faithfulness violation test for the analyzed explanation methods. The results are averaged on adopted models and datasets.
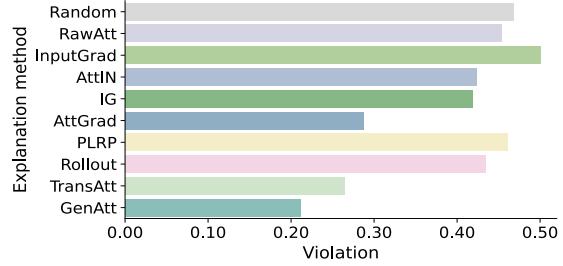
highly close. The correlation plot between the Sufficiency and Violation for the RawAtt can be a concrete example. As the faithfulness violation level increases, the scores on the Sufficiency, however, do not exhibit an increasing trend, *i.e.*, higher Sufficiency means worse faithfulness. Even in the large divergence between 0.0 and 1.0 violation levels, their distribution for the Sufficiency scores is still highly close. In this sense, we can see that existing metrics do not equip capabilities of examining the violation issue, which highlights the necessity of our violation test.

**Analysis of violation behavior.** To further verify that our test is responsible for the faithfulness violation detection, we compare the behavior of identified faithfulness violators and non-violators (See Figure 3). From the results, we can observe that in non-violators, the explanation weights and the measured impact for each input tokens are highly consistent. For example, for raw attention maps, the input tokens with higher positive weights actually show large contribution effects to the model predictions. This behavior is in stark with the detected violators, where the input with the largest weights (in absolute values) tends to present the opposite effects as the polarity indications, *e.g.*, heavy contribution effects are assigned with the most negative weights in Attention ⊙ Gradient. Such divergence in their behaviors demonstrates the effectiveness of our proposed test for the faithfulness violation examination.

In addition, it can be further observed that the relatively

lower explanation weights often exhibit the consistent polarity with the input impact, *e.g.*, inputs with rank-2 positive weights mostly impose contribution effects. This may be the reason why existing explanation metrics perform poorly in violation detection – their calculation is based on the test of a large proportion of input tokens, such that the polarity violation in the forefront weight can be easily concealed.

**5.2. Evaluation results for explanation methods (RQ2)**

For a complete faithfulness evaluation and comparison, we further examine explanation methods with all faithfulness evaluation metrics on all utilized models and datasets. Due to space constraints, we only provide visualized results for 3 datasets (QQP, VQA 2.0, GQA) in Figure 5. In particular, the average of violation test measures for explanation methods is illustrated in Figure 4. We provide the main observations in the following:

Firstly, we can observe that all tested explanation methods unexpectedly suffer from the faithfulness violation issue, *i.e.*, violation ratio > 0. While in some cases many methods hardly exhibit the violation behavior, they can still perform worse on other architectures or datasets (See Figure 5). For example, the explanation method IG on the BUTD model only shows 4% violation ratio; whereas the 23% violation ratio can be drawn from the CNN+TanhAtt model.

Secondly, among the tested attention-based explanation methods, RawAtt and AttIN perform the worst on the violation test. We consider this is due to the fact that their explanation weights are only with one polarity, *i.e.*, positive
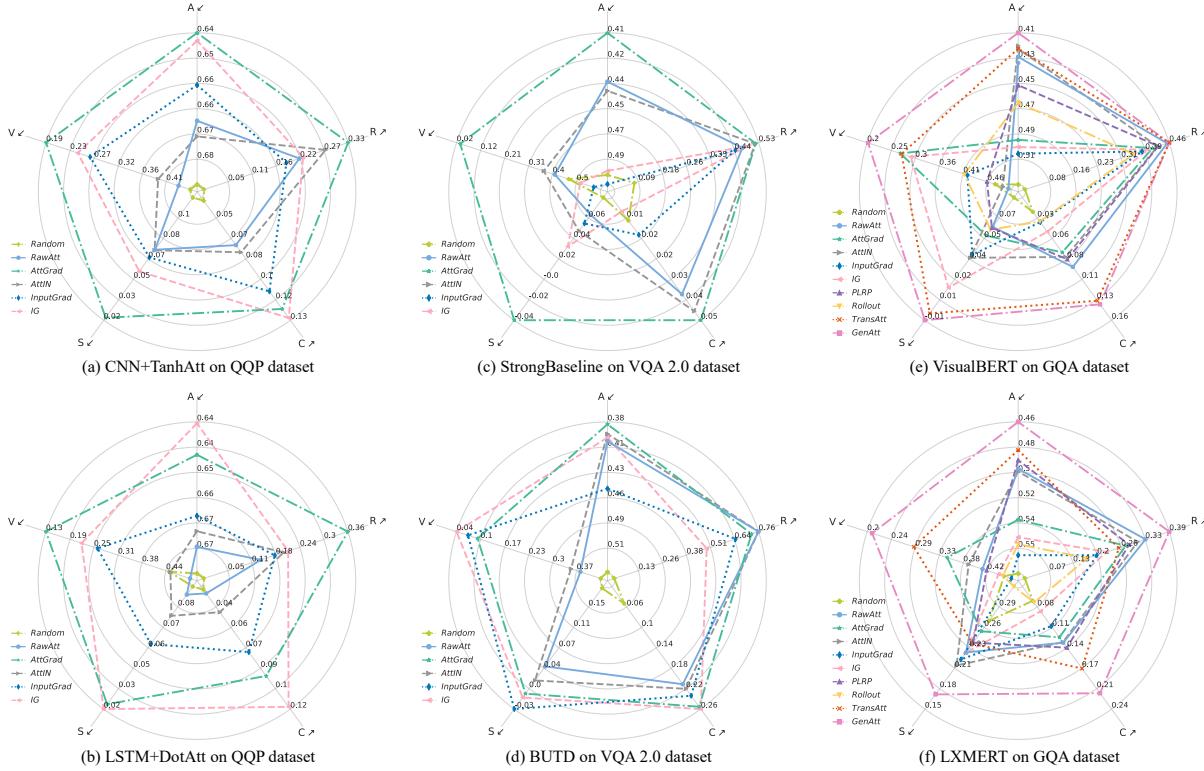
*Figure 5.* Faithfulness diagnosis for explanation methods over different architectures, on the QQP, VQA 2.0, GQA datasets. The ✓ signs indicate the lower value of the corresponding metric is better, while ↗ denotes the higher value is better. Table 3 exhibits the denotations of metrics. Appendix D provides the plots for all models and datasets, and the average evaluation scores for all explanation methods.

(See Figure 3 for an example of RawAtt weights). In this situation, their single-polarity weights do not equip the capability to identify whether the feature impact is positive or negative, thereby leading to the high faithfulness violation (We provide more discussion in Section 5.3). Moreover, we can further observe that even these two methods are heavily hindered by the violation issue, they commonly do not show poor performance on the other evaluation metrics (See Figure 5). For instance, in BUTD and LXMERT models, while RawAtt and AttIN present severe violations (about 40%), they can still show promising performance on the AUCTP (A) and Sufficiency (S), which is even comparable or superior to the AttGrad (Note that AttGrad performs much better than them on the violation test). This finding helps us further deduce that the execution of violation test is imperative for a sanity faithfulness evaluation.

Lastly, through the results, the form of the Attention ⊙ Gradient performs the best almost on all architectures and datasets. In detail, for the general attention models (*e.g.*, CNN+TanhAtt or StrongBaseline), the AttGrad method often shows the best performance on all metrics. Likewise, for the transformer-based models, the another explanation approach GenAtt behaves better.

## 5.3. Factors affecting faithfulness violation issue (RQ3)

The previous section shows that the existing explanation methods mostly violate the faithfulness property in terms of polarity consistency across models and datasets, while the gradient-based explanations can alleviate it to some extent. Herein, we take a step further to study the potential factors leading to the faithfulness violation issue.

**Do deeper layers hinder the explanation faithfulness?** As shown before in Figure 4 and Figure 5, the explanation methods which perform well on the general attention models (*e.g.*, CNN+DotAtt) often exhibit a worse faithfulness on the transformer-based models, *e.g.*, an increase of 14% violation ratio in AttGrad. In the light of the prior studies (Shrikumar et al., 2017; Sundararajan et al., 2017; Sixt et al., 2020), we hypothesize that the divergence in number of stacked deep layers *w.r.t* these two architectures may be the reason. To this end, we study the faithfulness of explanation methods when the number of deep layers increases on the Yelp and SST datasets (See Figure 6). Through the results, a positive correlation between the layer numbers and the faithfulness violation can be observed on almost all explanation methods over both datasets, which confirms our hypothesis in the negative effects of deep layers. Besides, as
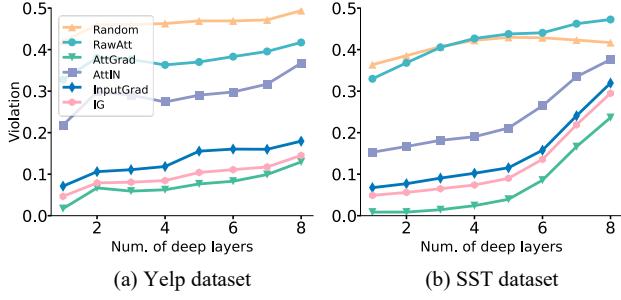
(a) Yelp dataset      (b) SST dataset

*Figure 6.* Faithfulness violation in relation to the number of non-linear deep layers. The results for Yelp and SST datasets are averaged on the four utilized attention models: CNN+TanhAtt, LSTM+TanhAtt, CNN+DotAtt, and LSTM+DotAtt.

model architecture grows in complexity and learning functions, this experiment further poses an imperative demand in developing robust explanation methods.

**Is it crucial to equip capability of identifying impact polarity?** In our previous findings, we show that AttIN and RawAtt are weak in identifying impact polarity because of their invariably positive weights. To verify the reason, we compare the explanation faithfulness of Attention ⊙ Gradient with the raw attention. As illustrated in Table 4, we ablate the formulation of Attention ⊙ Gradient ($\alpha \odot \nabla\alpha$) into two forms: 1) remaining only the sign of gradients ($\alpha \odot \text{sign}(\nabla\alpha)$), and 2) remaining merely the absolute value of gradients ($\alpha \odot |\nabla\alpha|$). By comparison with the raw attention explanations, we can observe that even with the signs of gradients only, the faithfulness still shows a large improvement over the raw attention. In addition, compared to the intact gradient information, the violation issue would also become more severe if we merely incorporate the absolute value of gradients. These observations spotlight the importance of capability to identify the impact polarity.

Furthermore, we also provide an example in Figure 7 to illustrate the above point. As can be seen, the given example consists of two parts with opposite sentiments (*i.e.*, positive from "good service" and negative from "horrible garbage"), which impose the contrary impacts on the model prediction. Suffered from the weakness of single-polarity weights, however, AttIN and RawAtt only can assign the same polarity for these two parts of words, which is therefore highly misleading and heavily violate the faithfulness.

## 6. Discussions and Conclusions

In this paper, we have presented the faithfulness violation view for efficient comparison of attention-based explanation methods. We have demonstrated that the polarity consistency defined with concrete mathematical formulation provides the reasonable evaluation criterion on the explainability of attention models. To facilitate the understanding of the

*Table 4.* Faithfulness violation ratios for ablations of attention-based explanation methods. $\alpha$ denotes the raw attention weights.

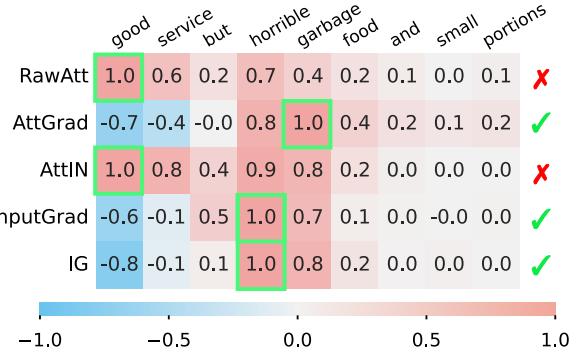| Method | Yelp | AgNews | VQA 2.0 |
|---|---|---|---|
| $\alpha$ | 0.31 | 0.28 | 0.40 |
| $\alpha \odot \nabla\alpha$ | **0.02** | **0.03** | **0.06** |
| $\alpha \odot \|\nabla\alpha\|$ | 0.15 | 0.07 | 0.25 |
| $\alpha \odot \text{sign}(\nabla\alpha)$ | 0.16 | 0.18 | 0.27 |



*Figure 7.* Normalized explanations for an example from the sentiment analysis task over the Yelp dataset. *Negative* sentiment is the model prediction for this example. Red boxes represent that indications of explanation weights are contribution effects to the model prediction. Blue boxes denote the indications are suppression effects. ✓: Non-violation. ✗: Violation.

faithfulness violation test, we have shown that our method is essential for the faithfulness evaluation in comparison with existing metrics (Section 5.1). Despite remarkable progress has been made in attention modeling, our results unexpectedly demonstrate that most tested explanation methods are hindered by the faithfulness violation issue in terms of the polarity consistency (Section 5.2). Nevertheless, this does not necessarily mean that the attention-based explanations examined are not meaningful, though the research has not yet reached the level of maturity. To gain useful insights in designing better models and explanations, we have studied the factors affecting the violation issue (Section 5.3), which suggests that the complexity in model architectures and the capability of identifying polarity can be the dominant factors. These impressive results highlight not only the current state of attention-based explanations, but also the challenges in designing and benchmarking future explanation methods.

Given the diversity of design philosophies, it is inevitably non-trivial to generalize across model architectures and tasks. Even though, the capability of identifying impact polarity should not be ignored in designing an explanation method. Without such capability, the explanations cannot tell us whether an input component would contribute to the model prediction or not, which can easily mislead us in understanding the rationales behind model decisions. While

incorporating the gradient-based explanation is not a flawless approach, our findings reveal that it does provide some clues and helps in polarity identification.

We acknowledge that the confidence perturbation can be the limitation in our faithfulness violation test, thus depressing the evaluation results to some extent. However, we would like to highlight again that our experiments are conducted by using several replacement functions, which are verified as promising solutions for OOD problem in perturbation. In addition, while there can be model noises (*e.g.*, biases) affecting violation identification, such errors are indeed unavoidable in current explanation evaluations, and our experimental results (*e.g.*, Figure 3) have also shown the potential of our tested method in faithfulness examination.

It is worth mentioning that the attention-based explanation methods discussed so far are constrained in an ideal narrow scope, facilitating the in-depth and focused discussion. However, there is a wider variety of applications and an enormous demand in attention model explanations. How to adapt the faithfulness violation to quantify the quality of explanation methods has great potential in improving model explainability. Lastly, we hope that we have provided further motivation to the community to consider the faithfulness violation as the design criteria in explanation methods and explainable attention models. We also envision that this is the most considerable benefit eventually lies.

## Acknowledgements

## References

Abnar, S. and Zuidema, W. H. Quantifying attention flow in transformers. In *ACL*, pp. 4190–4197. Association for Computational Linguistics, 2020.

Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I. J., Hardt, M., and Kim, B. Sanity checks for saliency maps. In *NeurIPS*, pp. 9525–9536. Curran Associates, Inc., 2018.

Ancona, M., Ceolini, E., Öztireli, C., and Gross, M. Towards better understanding of gradient-based attribution methods for deep neural networks. In *ICLR*, 2018.

Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., and Zhang, L. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*, pp. 6077–6086. IEEE, 2018.

Atanasova, P., Simonsen, J. G., Lioma, C., and Augenstein, I. A diagnostic study of explainability techniques for text classification. In *EMNLP* , pp. 3256–3274. Association for Computational Linguistics, 2020.

Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. In *ICLR*, 2015.

Bai, B., Liang, J., Zhang, G., Li, H., Bai, K., and Wang, F. Why attentions may not be interpretable? In *KDD*, pp. 25–34. ACM, 2021.

Binder, A., Montavon, G., Lapuschkin, S., Müller, K., and Samek, W. Layer-wise relevance propagation for neural networks with local renormalization layers. In *ICANN* , pp. 63–71. Springer, 2016.

Bowman, S. R., Angeli, G., Potts, C., and Manning, C. D. A large annotated corpus for learning natural language inference. In *EMNLP*, pp. 632–642. The Association for Computational Linguistics, 2015.

Brunner, G., Liu, Y., Pascual, D., Richter, O., Ciaramita, M., and Wattenhofer, R. On identifiability in transformers. In *ICLR*, 2020.

Chefer, H., Gur, S., and Wolf, L. Generic attention-model explainability for interpreting bi-modal and encoder-decoder transformers. In *ICCV*, pp. 397–406, 2021a.

Chefer, H., Gur, S., and Wolf, L. Transformer interpretability beyond attention visualization. In *CVPR*, pp. 782–791. IEEE, 2021b.

Chen, J., Song, L., Wainwright, M. J., and Jordan, M. I. L-shapley and c-shapley: Efficient model interpretation for structured data. In *ICLR*, 2019.

Chrysostomou, G. and Aletras, N. Improving the faithfulness of attention-based explanations with task-specific information for text classification. In *ACL/IJCNLP* , pp. 477–488. Association for Computational Linguistics, 2021.

Clark, K., Khandelwal, U., Levy, O., and Manning, C. D. What does BERT look at? an analysis of bert's attention. In *BlackboxNLP@ACL*, pp. 276–286. Association for Computational Linguistics, 2019.

DeYoung, J., Jain, S., Rajani, N. F., Lehman, E., Xiong, C., Socher, R., and Wallace, B. C. ERASER: A benchmark to evaluate rationalized NLP models. In *ACL*, pp. 4443–4458. Association for Computational Linguistics, 2020.

Ethayarajh, K. and Jurafsky, D. Attention flows are shapley value explanations. In *ACL/IJCNLP* , pp. 49–54. Association for Computational Linguistics, 2021.

Fong, R., Patrick, M., and Vedaldi, A. Understanding deep networks via extremal perturbations and smooth masks. In *ICCV*, pp. 2950–2958. IEEE, 2019.

Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., and Parikh, D. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, pp. 6325–6334. IEEE, 2017.

Hao, Y., Dong, L., Wei, F., and Xu, K. Self-attention attribution: Interpreting information interactions inside transformer. In *AAAI*, pp. 12963–12971. AAAI Press, 2021.

Hase, P., Xie, H., and Bansal, M. The out-of-distribution problem in explainability and search methods for feature importance explanations. In *NeurIPS*. Curran Associates, Inc., 2021.

He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *CVPR*, pp. 770–778. IEEE, 2016.

Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, 1997.

Hooker, S., Erhan, D., Kindermans, P., and Kim, B. A benchmark for interpretability methods in deep neural networks. In *NeurIPS*, pp. 9734–9745. Curran Associates, Inc., 2019.

Hudson, D. A. and Manning, C. D. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, pp. 6693–6702. IEEE, 2019.

Jacovi, A. and Goldberg, Y. Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness? In *ACL*, pp. 4198–4205. Association for Computational Linguistics, 2020.

Jain, S. and Wallace, B. C. Attention is not explanation. In *NAACL-HLT* , pp. 3543–3556. Association for Computational Linguistics, 2019.

Kazemi, V. and Elqursh, A. Show, ask, attend, and answer: A strong baseline for visual question answering. *arXiv preprint arXiv:1704.03162*, 2017.

Kobayashi, G., Kuribayashi, T., Yokoi, S., and Inui, K. Attention is not only a weight: Analyzing transformers with vector norms. In *EMNLP* , pp. 7057–7075. Association for Computational Linguistics, 2020.

LeCun, Y., Kavukcuoglu, K., and Farabet, C. Convolutional networks and applications in vision. In *ISCAS*, pp. 253–256. IEEE, 2010.

Li, J., Monroe, W., and Jurafsky, D. Understanding neural networks through representation erasure. *arXiv preprint arXiv:1612.08220*, 2016.

Li, L. H., Yatskar, M., Yin, D., Hsieh, C.-J., and Chang, K.-W. What does bert with vision look at? In *ACL*. Association for Computational Linguistics, 2020.

Lundberg, S. M. and Lee, S. A unified approach to interpreting model predictions. In *NeurIPS*, pp. 4765–4774. Curran Associates, Inc., 2017.

Luong, T., Pham, H., and Manning, C. D. Effective approaches to attention-based neural machine translation. In *EMNLP*, pp. 1412–1421. The Association for Computational Linguistics, 2015.

Michel, P., Levy, O., and Neubig, G. Are sixteen heads really better than one? In *NeurIPS*, pp. 14014–14024. Curran Associates, Inc., 2019.

Mnih, V., Heess, N., Graves, A., and Kavukcuoglu, K. Recurrent models of visual attention. In *NeurIPS*, pp. 2204–2212. Curran Associates, Inc., 2014.

Mohankumar, A. K., Nema, P., Narasimhan, S., Khapra, M. M., Srinivasan, B. V., and Ravindran, B. Towards transparent and explainable attention models. In *ACL*, pp. 4206–4216. Association for Computational Linguistics, 2020.

Moradi, P., Kambhatla, N., and Sarkar, A. Measuring and improving faithfulness of attention in neural machine translation. In *EACL*, pp. 2791–2802. Association for Computational Linguistics, 2021.

Pennington, J., Socher, R., and Manning, C. Glove: Global vectors for word representation. In *ENMLP*, pp. 1532–1543. Association for Computational Linguistics, 2014.

Ren, S., He, K., Girshick, R., and Sun, J. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, pp. 91–99. Curran Associates, Inc., 2015.

Rensink, R. A. The dynamic representation of scenes. *Visual Cognition*, 7(1-3):17–42, 2000. doi: 10.1080/135062800394667.

Ribeiro, M. T., Singh, S., and Guestrin, C. "why should I trust you?": Explaining the predictions of any classifier. In *KDD*, pp. 1135–1144. ACM, 2016.

Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. Grad-cam: Visual explanations from deep networks via gradient-based localization. *Int. J. Comput. Vis.*, 128(2):336–359, 2020.

Serrano, S. and Smith, N. A. Is attention interpretable? In *ACL*, pp. 2931–2951. Association for Computational Linguistics, 2019.

Shrikumar, A., Greenside, P., and Kundaje, A. Learning important features through propagating activation differences. In *ICML*, pp. 3145–3153. PMLR, 2017.

Sixt, L., Granz, M., and Landgraf, T. When explanations lie: Why many modified BP attributions fail. In *ICML*, pp. 9046–9057. PMLR, 2020.

Smilkov, D., Thorat, N., Kim, B., Viégas, F. B., and Wattenberg, M. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.

Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A. Y., and Potts, C. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*, pp. 1631–1642. Association for Computational Linguistics, 2013.

Srinivas, S. and Fleuret, F. Rethinking the role of gradient-based attribution methods for model interpretability. In *ICLR*, 2021.

Sundararajan, M., Taly, A., and Yan, Q. Axiomatic attribution for deep networks. In *ICML*, pp. 3319–3328. PMLR, 2017.

Swayamdipta, S., Schwartz, R., Lourie, N., Wang, Y., Hajishirzi, H., Smith, N. A., and Choi, Y. Dataset cartography: Mapping and diagnosing datasets with training dynamics. In *EMNLP*, pp. 9275–9293. Association for Computational Linguistics, 2020.

Tan, H. and Bansal, M. LXMERT: Learning cross-modality encoder representations from transformers. In *EMNLP/IJCNLP*, pp. 5099–5110. Association for Computational Linguistics, 2019.

Thorne, J., Vlachos, A., Christodoulopoulos, C., and Mittal, A. Generating token-level explanations for natural language inference. In *NAACL-HLT*, pp. 963–969. Association for Computational Linguistics, 2019.

Tomsett, R., Harborne, D., Chakraborty, S., Gurram, P., and Preece, A. D. Sanity checks for saliency metrics. In *AAAI*, pp. 6021–6029. AAAI Press, 2020.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. In *NeurIPS*, pp. 5998–6008. Curran Associates, Inc., 2017.

Voita, E., Talbot, D., Moiseev, F., Sennrich, R., and Titov, I. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *ACL*, pp.

5797–5808. Association for Computational Linguistics, 2019.

Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *ICLR*, 2019.

Weston, J., Bordes, A., Chopra, S., and Mikolov, T. Towards ai-complete question answering: A set of prerequisite toy tasks. In *ICLR*, 2016.

Wiegreffe, S. and Pinter, Y. Attention is not not explanation. In *EMNLP/IJCNLP*, pp. 11–20. Association for Computational Linguistics, 2019.

Yang, Z., He, X., Gao, J., Deng, L., and Smola, A. J. Stacked attention networks for image question answering. In *CVPR*, pp. 21–29. IEEE, 2016.

Zeiler, M. D. and Fergus, R. Visualizing and understanding convolutional networks. In *ECCV*, pp. 818–833. Springer, 2014.

Zintgraf, L. M., Cohen, T. S., Adel, T., and Welling, M. Visualizing deep neural network decisions: Prediction difference analysis. In *ICLR*, 2017.

# A. Tasks and datasets

- **Sentiment Analysis**. We use Stanford Sentiment Treebank (SST) (Socher et al., 2013) and Yelp[2] datasets. In this task, the trained models are required to predict sentiment for each sentence from five levels (1 for most negative and 5 for most positive).

- **Topic Classification**. We utilize AGNews[3] and 20News[4] datasets, where the target is to predict categories (*e.g.*, business) for each article. In particular, AGNews contains 496,835 articles from 2000+ sources, while 20News consists of 20,000 samples.

- **Paraphrase Detection**. We adopt the Quora Question Paraphrase (QQP) dataset (Wang et al., 2019) following (Mohankumar et al., 2020). In this task, models are trained to determine the semantic equivalence of a pair of questions.

- **Natural Language Inference**. We utilize the SNLI (Bowman et al., 2015), which includes about 570K human-written English sentence pairs. In this task, models are trained for identifying textual entailment within sentence pairs.

- **Question Answering**. We make use of the bAbI-1 dataset (Weston et al., 2016) here. In this dataset, correctly answering each question requires a supporting fact from the context.

- **Visual Question Answering (VQA)**. We use VQA 2.0 (Goyal et al., 2017) and GQA (Hudson & Manning, 2019) datasets, which include about 410K and 950K training samples, respectively. In this task, models are required to understand both images and texts for correctly answering questions.

# B. Details of explanation methods

We mainly study the attention model explainability from three groups of explanation methods – *generic attention-based*, *transformer-based*, and *gradient-based* attribution approaches. In detail, the following methods are utilized.

### B.1. Generic attention-based explanation methods

Starting from generic attention-based explanation methods, we employ them in all attention models.

**Inherent Attention Explanation (RawAtt)** (Jain & Wallace, 2019) is the direct output from attention mechanisms. In *AddAtt*, for example, it is calculated by $\alpha = \mathrm{softmax}(\mathbf{w}_1^T \tanh(\mathbf{W}_2 K + \mathbf{W}_3 Q))$, where K and Q are query and keys. $\mathbf{w}_1$, $\mathbf{W}_2$, and $\mathbf{W}_3$ are learned parameters.

**Attention $\odot$ Gradient (AttGrad)** (Serrano & Smith, 2019; Chrysostomou & Aletras, 2021; Hao et al., 2021) scales attention weights with the corresponding gradients, denoted as $\alpha \odot \nabla\alpha$, where $\nabla\alpha = \frac{\partial \hat{y}}{\partial \alpha}$.

**Attention $\odot$ InputNorm (AttIN)** (Kobayashi et al., 2020) incorporates another factors, input vectors, in the attention-based explanations. *InputNorm* corresponds to the norm of transformed input vectors, $||v(x)||$, in attention modules, where $v(\cdot)$ can be the value mappings in transformers.

### B.2. Transformer-based explanation methods

Regarding the exclusive architecture of transformers (*e.g.*, multi-head attention), we further adopt explanation methods for transformer-based models in the following.

**Partial Layer-wise Relevance Propagation (PLRP)** (Voita et al., 2019) computes the LRP-based scores (Binder et al., 2016) for each attention head in the last layer[5], and then average them across heads as explanations.

**Rollout** (Abnar & Zuidema, 2020) builds upon the assumption that attention weights can be combined linearly following paths along the pairwise attention graph.

**Transformer Attention Attribution (TransAtt)** (Chefer et al., 2021b) is another form of LRP-based explanation method. It calculates LRP scores for each attention heads, and then integrates such scores with attention gradients throughout the

---

[2]www.yelp.com/dataset_challenge

[3]http://www.di.unipi.it/ gulli/AG_corpus_of_news_articles.html

[4]http://qwone.com/ jason/20Newsgroups

[5]LRP is a method to compute contribution scores for neurons.

attention graph. We simplify it as $\nabla\alpha \odot \mathrm{R}^\alpha$, where $\mathrm{R}^\alpha$ is the LRP scores for the corresponding head.

**Generic Attention Attribution (GenAtt)** (Chefer et al., 2021a) provides general solution to explain different transformer architectures. GenAtt combines attention weights with corresponding gradients for all layers based on the devised adaptation strategies, and therefore can be deemed as another form of the $\alpha \odot \nabla\alpha$.

### B.3. Gradient-based attribution methods

Besides the above methods, we also select two representatives from gradient-based approaches for a comparison:

**Input $\odot$ Gradient (InputGrad)** (Shrikumar et al., 2017; Ancona et al., 2018) calculates the element-wise product between the input $x$ and the corresponding gradient: $x \odot \frac{\partial \hat{y}}{\partial x}$.

**Integrated Gradients (IG)** (Sundararajan et al., 2017) alleviates gradient saturation issue by aggregating input gradients along the path between a baseline input $(x')$ to the original input $(x)$, defined as $(x - x') \times \int_0^1 \frac{\partial f(x' + \gamma(x - x'))}{\partial x} d\gamma$.

## C. Details of faithfulness evaluation metrics

**AUC-TP** $\downarrow$ (Atanasova et al., 2020). The calculation of this metric is based on the AUC score of performance perturbation. Concretely, we first sequentially mask the top-0, 5, 10, 20, ..., 90% input tokens *w.r.t.* explanation weights, to obtain a set of model performance $S = \{s^0, s^5, s^{10}, s^{20}, ..., s^{90}\}$. Then based on the mask levels, we calculate the area under threshold-performance curve, defined as $\mathrm{AUC}([0, 5\%, 10\%, 20\%, ..., 90\%], S)$.

**Sufficiency** $\downarrow$ (DeYoung et al., 2020) measures whether important features identified by the explanation method are adequate to remain confidence on the original predictions. Formally,

$$\text{Sufficiency} = \frac{1}{B} \sum_{k \in B} f(x)_{\hat{y}} - f(x_{:k\%})_{\hat{y}}, \tag{3}$$

where $B = \{5, 10, 20, 50\}$, representing sparsity levels of explanation weights in order of decreasing explanation weights.

**Comprehensiveness** $\uparrow$ (DeYoung et al., 2020) evaluates if the features assigned lower weights are unnecessary for the predictions. Formally,

$$\text{Comprehensiveness} = \frac{1}{B} \sum_{k \in B} f(x)_{\hat{y}} - f(x \backslash x_{:k\%})_{\hat{y}}, \tag{4}$$

where $B = \{5, 10, 20, 50\}$, denoting sparsity levels of explanation weights in order of decreasing explanation weights.

**Rank Correlation (RC)** $\uparrow$ (Jain & Wallace, 2019; Selvaraju et al., 2020) measures the monotonic correlation between the explanation weights and the feature importance. The calculation of this metric is based on the Spearman's Rank Correlation Coefficient, *i.e.*, $\rho(\cdot)$. Formally,

$$\text{Rank-Correlation} = \rho(\hat{e}, p), \tag{5}$$

where $\hat{e}$ denotes the absolute explanation weights in the descending order, and $p$ represents the corresponding measured feature importance, *i.e.*, $p = \left[ \sum_i |f(x)_{y_i} - f(x \backslash x_{:1})_{y_i}|, \sum_i |f(x)_{y_i} - f(x \backslash x_{:2})_{y_i}|, ..., \sum_i |f(x)_{y_i} - f(x \backslash x_{:N})_{y_i}| \right]$.

## D. Sanity faithfulness evaluation for explanation methods

*Table 5.* Average faithfulness evaluation scores for explanation methods on all utilized datasets and model architectures.

| Method | Faithfulness Properties | | | | |
|---|---|---|---|---|---|
| | AUCTP↓ | Violation↓ | Suf.↓ | Comp.↑ | RC↑ |
| *Explanation for General Attention Models* | | | | | |
| Random | 0.570 | 0.425 | 0.226 | 0.060 | -0.030 |
| RawAtt | 0.422 | 0.309 | 0.105 | 0.248 | 0.411 |
| InputGrad | 0.386 | 0.121 | 0.033 | 0.310 | 0.437 |
| AttIN | 0.402 | 0.244 | 0.089 | 0.257 | 0.478 |
| IG | **0.345** | 0.095 | **-0.004** | **0.370** | 0.470 |
| AttGrad | 0.352 | **0.069** | 0.018 | 0.359 | **0.535** |
| *Explanation for Transformer-based models* | | | | | |
| Random | 0.581 | 0.468 | 0.189 | 0.028 | 0.000 |
| RawAtt | 0.477 | 0.454 | 0.150 | 0.136 | 0.360 |
| InputGrad | 0.560 | 0.501 | 0.159 | 0.090 | 0.321 |
| AttIN | 0.472 | 0.423 | 0.146 | 0.146 | 0.136 |
| IG | 0.551 | 0.419 | 0.139 | 0.120 | 0.318 |
| AttGrad | 0.538 | 0.287 | 0.145 | 0.127 | 0.131 |
| PLRP | 0.471 | 0.461 | 0.153 | 0.148 | 0.301 |
| Rollout | 0.533 | 0.434 | 0.190 | 0.073 | 0.105 |
| TransAtt | 0.458 | 0.265 | 0.123 | 0.194 | 0.299 |
| GenAtt | **0.451** | **0.212** | **0.093** | **0.318** | **0.434** |

(a) VisualBERT on VQA 2.0 dataset

(c) VisualBERT on GQA dataset

(b) LXMERT on VQA 2.0 dataset
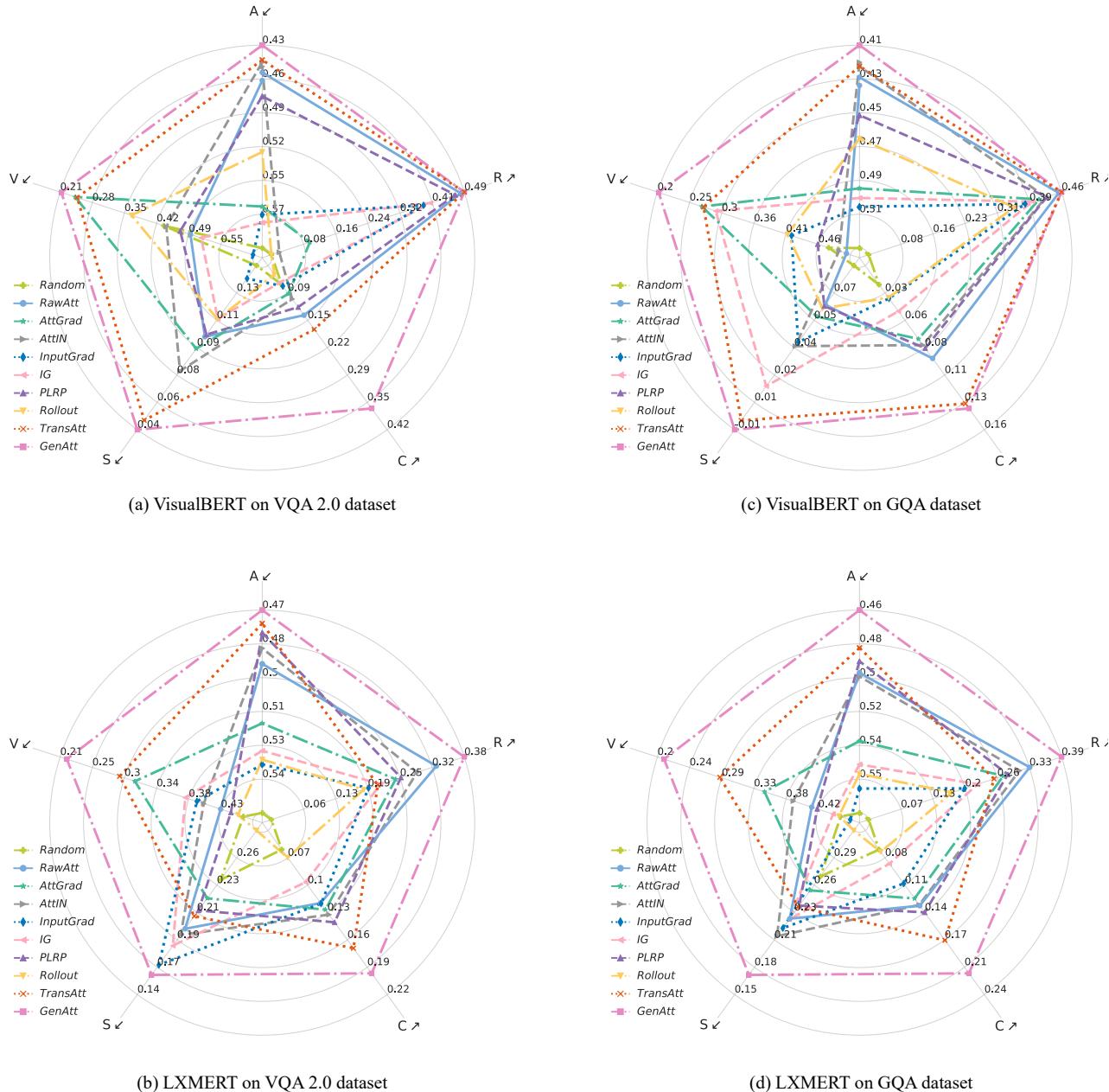
(d) LXMERT on GQA dataset

*Figure 8.* Faithfulness diagnosis for all explanation methods over different attention-based architectures on VQA 2.0 and GQA datasets.
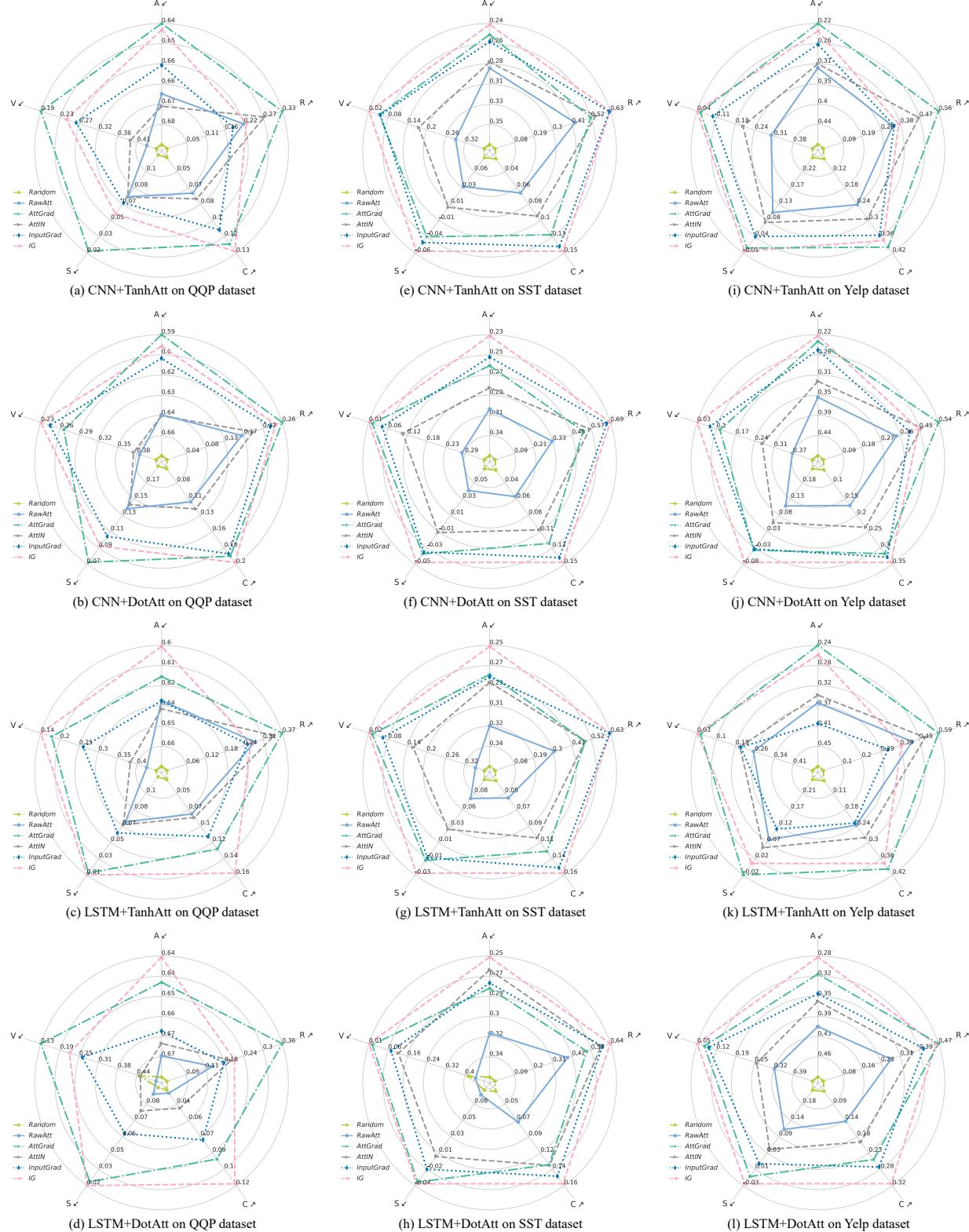
*Figure 9.* Faithfulness diagnosis for all explanation methods over different attention-based architectures on QQP, SST, and Yelp datasets.

Figure 10. Faithfulness diagnosis for all explanation methods over different attention-based architectures on the SNLI, bAbI-1, and AgNews datasets.
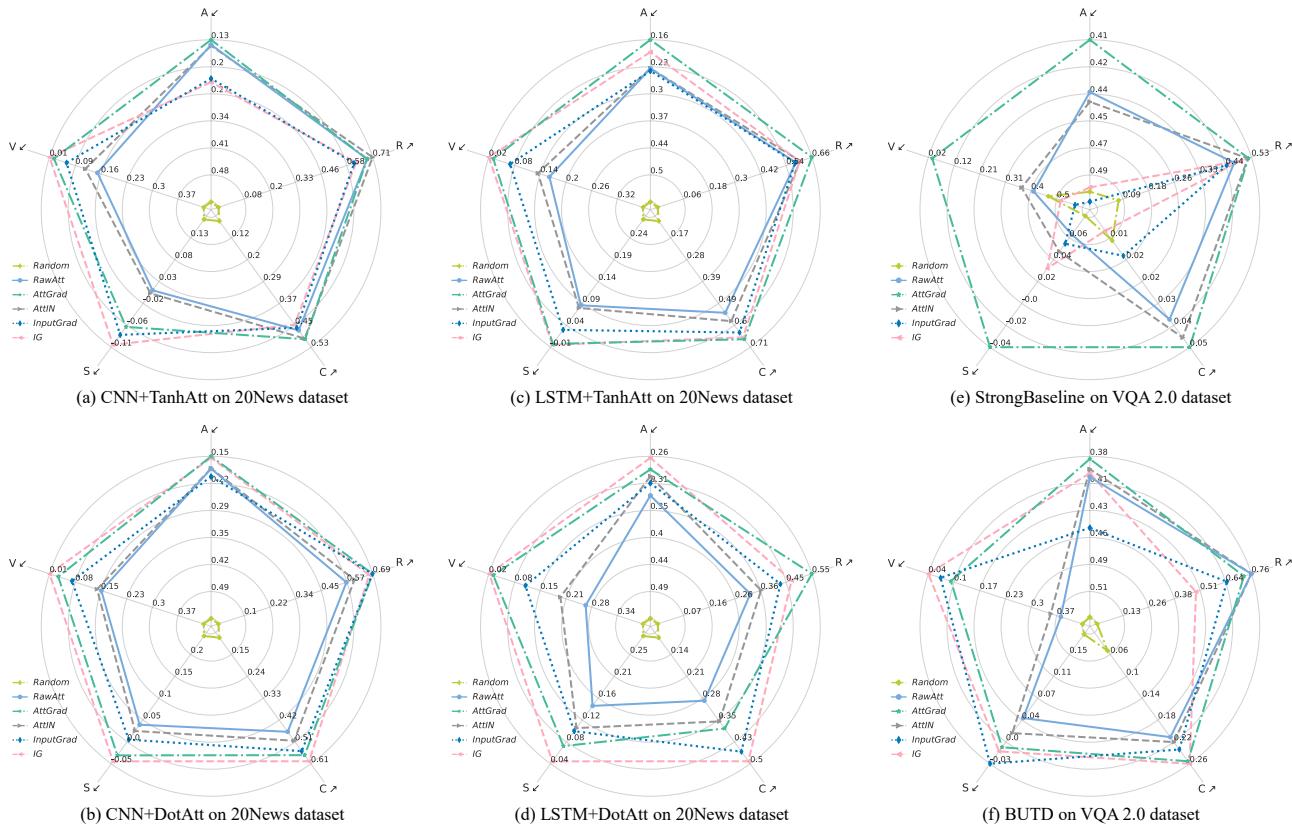
*Figure 11.* Faithfulness diagnosis for all explanation methods over different attention-based architectures on 20News and VQA 2.0 datasets.