Borislav Kolev Marinov B.K.Marinov-1@student.tudelft.nl

Evaluating Faithfulness of LLM Generated Explanations for Claims: Are Current Metrics Effective?

Analysing the Capabilities of Automatic Metrics to Represent the Difference in Faithfulness Between Explanations



01. Introduction

- Data Access → Claims and evidence flow into LLMs
- ● LLM → Produces outputs and explanations • Can sound convincing but be misleading
- \checkmark Experts ask \rightarrow "Why should we trust this explanation?" → evaluation required
- 🔹 🎯 Our Goal
 - Evaluate if existing metrics (G-Eval, FactCC, UniEval, QAGs) truly capture faithfulness in factchecking explanations

02. Research Question Terminology

How well do current evaluation metrics reflect the faithful ness of LLM-generated factchecking explana tions compared to journalist-written ones?

II Faithfulness Metrics

- 🗣 Claim Statement to be verified
- 듣 Evidence Factual context for claim
- 6 Label Verdict (True / Half-True / False)
- Q Faithfulness Explanation aligns only with evidence (no extra or false info)
- G-Eval LLM grades explanations using reasoning prompts with form-filling and Chain-of-Thought
- *FactCC* Checks consistency with evidence using BERT
- ? QAGs Converts explanation into questions, answers them with evidence
- VIIEval LLM uses boolean QA prompts to judge factual correctness
- Each metric gives a score between 0 (unfaithful) and 1 (faithful)



05. Analysis

- II Metric Correlation with Semantic Similarity
 - Weak correlations overall ($\rho \le 0.23$)
 - ॰ 🔽 FactCC & QAGs: Best (though limited) alignment 🖞 with expert explanations
 - 1 UniEval & G-Eval: Low or even negative correlations
- 🧠 Score Bias
 - All metrics assign slightly higher scores to incorrect predictions
 - 🙅 G-Eval shows the strongest bias favoring LLM outputs
- Perturbation Sensitivity
 - *X Unrelated Sentences: FactCC & UniEval show clear score drops
 - 🔺 🖹 Unsupported Sentences: Metrics often fail to detect hallucinations, though some score reduction observed

04. Results predictions, but not consistently Score Distribution: All Metrics for Generated Explanatio

Metric Correlations: Metrics favor LLMs – Normal distribution of scores for expert, Normal distribution skewed towards higher scores

Targeted Tests: Some metrics detect hallucinations (FactCC, UniEval) XX; others (G-Eval) fail to penalize unfaithful content consistently, while QAGs is robust to noise, but low in scores.



06. Conclusion

- Existing metrics are unreliable and inconsistent
- Different metrics perform better under the experiments, but no metric consistently reflects faithfulness
- Have some positive aspects that show latent potential
- Further Work
 - 💡 Use more capable LLMs (e.g., GPT-4) for generation and evaluation.
 - 🛠 Fine-tune metrics on the actual dataset (e.g., PolitiFact QuanTemp) to improve alignment.
 - Fest new perturbations insert noise or hallucinations into the middle of explanations.

Pradeep Murukannaiah, Shubhalaxmi Mukherjee



Accuracy to Fathfulness Scores: All metrics show slightly higher average scores for incorrect labels. Medians are higher for correct





