

Methods to Measure and Reduce Misinformation in Large Language Models

Machado

New Technologies in Data Analysis

HES-SO Master

Lausanne, Switzerland

.leitemac@hes-so.ch

Abstract—Large language models revolutionized natural language processing and the way humans interact with data by enabling advanced applications such as text generation, translation, and summarization. Their potential is unfortunately facing significant challenges, particularly in the generation of misinformation, biases and hallucination in the responses. It impacts their reliability, especially in critical domains like education, medicine or politics where accuracy is the main goal to achieve in the information retrieval. This paper provides a comprehensive review of methodologies designed to detect and reduce misinformation in LLMs. It explores established techniques, including Natural Language Inference, Question-Generation and Answering or formal methods, alongside emerging solutions such as Retrieval-Augmented Generation. These approaches are evaluated based on their strengths and limitations, highlighting the need for hybrid frameworks that combine multiple strategies to enhance reliability. The study also addresses the challenges of benchmarking and the rapid evolution of LLM technologies, proposing avenues for standardization and improvement. By focusing on integrated solutions and future research directions, this review aims to provide valuable insights for researchers and practitioners seeking to align LLMs more closely with human values and ethical standards. Through these advancements, the reliability and safety of LLM applications can be significantly improved, adopting them in diverse and sensitive fields.

Index Terms—LLM, Misinformation, RAG, Fake-news, RLHF

I. INTRODUCTION

Large Language Models such as GPT-4 have revolutionized natural language processing (NLP) methods by enabling capacities like coherent text generation, translation or summarization. It changed and probably will continue to change the way humans interact with computers and technologies and could arguably be described as one of the most important technology revolution in the past years. Despite their advancements, LLMs face significant challenges. Particularly, the generation of erroneous or incoherent outputs can be dangerous and limit their applicability in critical domains such as medicine, law or education. In these fields, accuracy and reliability are the first priority.

For instance, in the paper “Walking a Tightrope – Evaluating Large Language Models in High-Risk Domains” [1], authors describe a case where one model suggested cleaning a wound after a dog bite as adequate treatment, while the correct guidance should have been to seek medical attention for a possible rabies vaccination. This error highlights the risks of

using LLMs in medical contexts without robust oversight, as such misinformation could lead to serious consequences.

Hallucinations, including contradictions, factual or contextual inaccuracies or nonsensical outputs, represent a substantial risk to the broad adoption of these models. These errors often come from intrinsic limitations, such as biased or outdated source data, overfitting or superficial understanding or complex concepts. Furthermore, biases and undesirable behaviors embedded in the models raise ethical concerns. These limitations highlight the pressing need for robust methodologies to evaluate and mitigate misinformation generated by LLMs.

This paper investigates the challenges posed by LLMs in terms of misinformation and proposes a comprehensive review of existing techniques to address these issues. Drawing on both established practices and emerging technologies, it evaluates detection and mitigation methods, including Natural Language Inference, Question-Answering, and formal verification approaches. Additionally, it explores advanced solutions such as Retrieval-Augmented Generation (RAG) and activation-based truth prediction systems like SAPLMA.

The primary objective of this review paper is to offer a structured analysis of the technological challenges associated with misinformation and its evaluation in LLMs while identifying potential avenues for improvement.

By addressing these aspects, the study aims to provide precious insights for researchers and practitioners seeking to enhance the reliability and alignment of LLMs with human values.

The following study is structured as follows. Section II provides a comprehensive background on the foundational concepts of Large Language Models. In Section III, we explore various methodologies for detecting misinformation. Section IV delves into advanced techniques for reducing misinformation and concludes with an evaluation of these methods, highlighting their strengths and limitations through a comparative analysis. Finally, Section V and VI ends the paper with a discussion on future research questions and aspects of the ways to deal with misinformation in LLMs.

II. BACKGROUND

To understand the challenges and solutions addressed in this paper, it is crucial to define the key concepts and technologies used in the application and functioning of large

