# Hallucination Analysis in Large Language Models for Arabic Question Answering

Da'ad Albahdal[1] , Raghad Alawaji[1] , Murad A. Rassam[1,2]*

[1] Department of Information Technology, College of Computer, Qassim University, Buraydah 51452, Saudi Arabia
[2] Faculty of Engineering and Information Technology, Taiz University, Taiz 6803, Yemen

Corresponding Author Email: M.Qasem@qu.edu.sa

**ABSTRACT**

Large Language Models (LLMs) have demonstrated strong capabilities in generative and knowledge-intensive tasks such as question answering. However, these models may produce hallucinated responses, generating information that appears plausible but is factually incorrect or unsupported by evidence. While hallucinations in LLMs have been widely studied in English, limited attention has been given to their behavior in other languages, particularly Arabic. This study investigates hallucination phenomena in LLMs within the context of Arabic question answering. Two experimental settings are designed to evaluate the tendency of different models to generate hallucinated responses. In the first experiment, five widely used LLMs—Gemini, Claude 3.5, ChatGPT-4.0, ChatGPT-3.5, and Jais—are evaluated using hallucination-triggering questions in both Arabic and English. The results indicate that Gemini and Claude demonstrate stronger robustness against induced hallucinations, while ChatGPT-4.0 maintains consistent performance across both languages. The second experiment adopts a sampling-based framework to detect hallucinations in a black-box setting without relying on external knowledge sources or annotated references. The SelfCheckGPT framework is employed to evaluate GPT-4o-mini on open-domain Arabic question answering tasks. The model achieves an average hallucination score of 19.63% based on automatic sentence similarity evaluation. The findings highlight the challenges of hallucination in Arabic language processing and emphasize the need for more robust evaluation and mitigation strategies for LLMs in low-resource languages.

## 1. INTRODUCTION

Recently, Natural Language Processing (NLP) has gained significant popularity and is continually evolving. As a crucial link between humans and computers, NLP enables machines to comprehend, translate, and generate human language [1]. This evolution has contributed to the emergence of Large Language Models (LLMs) like BARD (Gemini Now) and ChatGPT, which enhance conventional chatbot capabilities through Deep Learning (DL) and transformer architectures [2]. LLMs find applications across various domains, including search engines, translation, medical diagnostics [3], code generation [4], higher education [3, 5], finance [6], and sport [7]. However, despite their impressive performance, a key challenge known as "hallucinations" hinders their practical reliability and precision as sources of information. In the field of computing, hallucinations refer to instances in which LLMs produce responses that sound convincing but are actually incorrect or nonsensical [8, 9]. Simply put, these hallucinations happen when the model creates information that seems believable yet isn't grounded in reality. This can include fabricating details, making up events, or straying beyond what it has learned from its training data, leading to inaccurate outputs. These responses may appear authentic but lack any basis in real-world facts [10]. Tackling this issue is essential for improving the reliability and trustworthiness of AI-generated content. Hallucination poses a significant challenge in medicine [3, 11], business decision-making [12], and other fields that require highly accurate measurements [13]. Looking back to the root causes of hallucination LLMs, researchers found that models deliver statistically based outcomes due to limited real-world knowledge, bias, or deceptive training data. This suggests a potential lack of understanding of input [14]. Furthermore, Hallucinations in LLMs affect their significant economic potential and rapid adoption. Recent research conducted by McKinsey in June 2023 highlights the economic growth and productivity gains enabled by AI generative applications. According to their estimates, this technology has the potential to annually generate value for the global economy ranging from $2.6 trillion to $4.4 trillion [15]. The introduction of ChatGPT by OpenAI in late 2022 garnered widespread attention for chatbot technology among the public. Within just five days of its launch, it amassed one million users. Notably, a milestone was achieved in January 2023, with the platform attracting over 100 million users [16].

Moreover, the growing adoption of LLMs is reaching diverse communities, including Arabic speakers, representing a significant segment of internet users. Arabic is the fourth most spoken language in the world, with over 480 million

native speakers [17, 18]. Approximately 67% of the Arab world's population are internet users [19], increasing demand for language processing technology. Evaluating LLMs in the context of the Arabic language has both practical applications and significant commercial relevance. Hence, developing LLMs that can understand and generate Arabic is vital for inclusive AI advancements. However, challenges such as a lack of high-quality Arabic datasets and dialectal variations hinder these models' effectiveness [20]. Consequently, hallucinations in AI outputs lead to inaccuracies and risk spreading misinformation among Arabic users, undermining the reliability of AI applications.

Current research works in the field of LLM hallucinations are examining the underlying factors and exploring methods to mitigate them. Detecting misinformation generated by LLMs is more challenging than detecting that generated by humans [21]. This highlights the importance of understanding the core of the problem and how different LLM models behave. Most LLM models are developed for Latin Languages, making it challenging to handle non-Latin environments [22]. Moreover, the size of training data plays a critical role in effective models where many non-Latin Languages have smaller sizes compared to Latin Languages [23-28]. To our knowledge, no previous academic studies have been conducted on the phenomenon of hallucinations of LLM in the Arabic language to date.

To address this gap, this research studies the phenomenon of hallucinations generated by LLMs in Arabic language, focusing on question-answering tasks with the primary aim of contributing to the continued advancement of the field. The main contributions of this study are as follows:

- To investigate hallucinations in Arabic language QA task in public prominent LLMs: Gemini, ChatGPT-4.0, ChatGPT-3.5, Claude 3.5, and Jais.
- To compare the tendency of hallucinations in the QA task in five prominent public LLMs for Arabic and English languages.
- To evaluate the hallucination rate of a black-box model adopting a sampling-based approach.
- To analyze the hallucination level in various generated responses.

The remaining sections of the paper are organized as follows: Section 2 presents the problem statement and research motivation. Section 3 provides background on NLP and Hallucination in LLMs. After that, Section 4 discusses recent studies on hallucination in LLMs across different tasks. Section 5 then introduces the research methodology. Following the presentation of the results and discussion in Section 6. Finally, Section 7 presents the conclusions and potential future work.

## 2. PROBLEM STATEMENT AND RESEARCH MOTIVATION

The present issue concerns Hallucination in LLMs, which occurs when these models produce text that appears credible but strays from factual truth or input faithfulness despite exhibiting fluency and grammatical correctness [23, 29]. Research has categorized and analyzed these hallucinations across various tasks, including machine translation, question answering, dialogue systems, summarization systems, knowledge graphs, and visual question-answering [30]. However, LLMs are increasingly being utilized by people who

require essential information, such as researchers [5], business sectors [12], students and teachers [31], and patients [32], due to their excellent performance. LLM hallucinations raise serious concerns about their reliability in real-world scenarios, as they significantly hinder practical use. This has prompted increased attention from researchers and practitioners who are focused on developing strategies to identify and address these hallucinations, aiming to enhance the trustworthiness and effectiveness of LLMs in practical applications. Moreover, most models are excessively trained in Latin Languages, making handling dialectal variations and morphological structures in non-English settings challenging since the effectiveness of model performance depends on the size of available training data [22, 23]. These factors emphasize the importance of examining LLM hallucination in language-specific domains to better understand and address the phenomenon effectively.

This study delves into the scope of LLMs hallucinations, narrowing its focus to hallucinations in natural language generation, specifically LLM transformer-based architecture with question-answering tasks in the Arabic language.

## 3. BACKGROUND

Natural Language Processing (NLP) is an AI field that processes speech and text to understand their syntactic, semantic, and emotional aspects [29]. Rule-based techniques were initially used [30, 31], followed by supervised and unsupervised approaches [32-34], and most recently, transformers bring a breakthrough in NLP tasks based on an encoder-decoder structure along with long-range dependency capture [35-37]. The development of LLMs has revolutionized NLP, achieving human-level performance on various language tasks, including medical diagnostics [3], code generation [4], higher education [3, 5], and Finance [6]. Despite the outstanding performance of these models, a fundamental issue called "hallucinations" has been preventing their practical usage as a reliable source of information. The subsequent sections outline Hallucination in LLMs, including its taxonomy, causes, and mitigation strategies.

### 3.1 Hallucination definition

In a computing context, the term hallucination is defined as " a plausible but false or misleading response generated by an artificial intelligence algorithm" [8]. The LLM Hallucinations occur when a model produces false or fabricated information, straying from factual knowledge and giving answers unsupported by the model's training data or the input provided [9]. This phenomenon can arise due to various factors related to the data, training processes, or inference methodologies, as discussed in subsection 3.3.

### 3.2 Classification of hallucinations

The categorization of LLM hallucinations remains an evolving field. While no single taxonomy reigns supreme, researchers have proposed diverse categories to classify them. This section reviews the different classification approaches. Figure 1 offers a detailed overview of the various classifications of hallucinations based on reviewed studies. Also, examples for each category are presented in Figure 2.
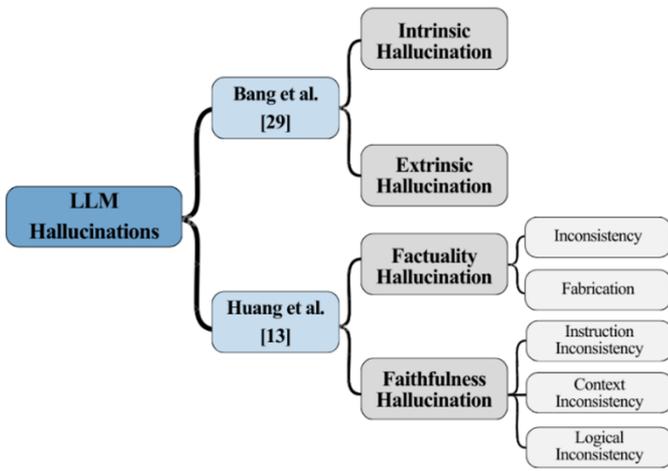
**Figure 1.** Overview of hallucination classification

natural language generation. The first category is Intrinsic hallucination which happens when an LLM produced output contradicts the original content. In contrast, the second category, extrinsic hallucination, arises when LLM produces outcomes that cannot be validated against the original content and lack support or contradiction within the original context.

Another study [14] introduced different types of Hallucination by categorizing hallucinations into two main groups. The first is factuality hallucination, which highlights the disparity between the outcomes produced and verified real-world facts. This category is divided into two subcategories based on verifiable sources. The first subcategory is factual inconsistency, which occurs when the AI provides outputs that include facts that can be confirmed but don't quite add up. The second subcategory is factual fabrication, which occurs when the AI presents information that cannot be verified and contradicts what is known to be true. This type can sometimes appear as either unverifiability or an overclaim of facts.

A study [30] has identified two forms of hallucinations in

| Category: Subcategory | Example | Explanation |
|---|---|---|
| Intrinsic Hallucination | • The provided data states that a city's population is 2.5 million, but the generated text reports it as 3 million. | The generated text includes information that goes against the data provided. |
| Extrinsic Hallucination | • The generated text claims that the city's population has grown by 10% in the last year, even though that information was not in the data provided. | The generated text includes extra information that isn't related to the data provided. |
| Factuality Hallucination: Inconsistency | • Instruction: What are Thomas Edison's main contributions to science and technology? • Response: Thomas Edison developed the first practical telephone and invented the light bulb. | The statement is factually inaccurate. In truth, Edison enhanced the light bulb by building upon the previous designs made by others, while Alexander Graham Bell was the one who invented the telephone. |
| Factuality Hallucination: Fabrication | • Instruction: What are the major environmental impacts of the construction of the Eiffel Tower? • Response: The construction of the Eiffel Tower in 1889 led to the extinction of the Parisian tiger, a species that played a crucial role in the region's ecosystem. Additionally, it is widely recognized as the event that sparked the global green architecture movement. | The response incorrectly claims that the construction of the Eiffel Tower led to the extinction of the fictional 'Parisian tiger' and initiated the global green architecture movement. There is no evidence of a 'Parisian tiger' existing, making this a fabricated claim. Furthermore, attributing green architecture solely to the Eiffel Tower oversimplifies its diverse origins. |
| Faithfulness Hallucination: Instruction Inconsistency | • Instruction: Translate the English question into Spanish: "What is the capital of France?" • Response: The capital of France is Paris | The user explicitly provided an instruction to translate the English question into Spanish. Instead of performing the translation as instructed, the LLM chose to answer the question. |
| Faithfulness Hallucination: Context Inconsistency | • Retrieved context: The Nile is a major north-flowing river in northeastern Africa. It's traditionally considered the longest river in the world. Originating in the Great Lakes region of central Africa, it flows through several countries before emptying into the Mediterranean Sea. • Instruction: Summarize the retrieval document above. • Response: The Nile originates from the mountain ranges of central Africa and flows into the Mediterranean | The summary incorrectly states that the Nile originates from the "mountain ranges of central Africa", while the paragraph mentions that the Nile originates in the Great Lakes region of central Africa. |
| Faithfulness Hallucination: Logical Inconsistency | • Instruction: Solve the equation 2x + 3 = 11 step by step. • Response: Step 1: Subtract 3 from both sides to get: 2x = 8. Step 2: Divide both sides by 2 to get: x = 3. So the result of the equation is x=3. | The first step correctly isolates the term with x to get 2x = 8. However, the next step inaccurately divides 8 by 2 to yield a result of x = 3, which is inconsistent with the earlier reasoning. |

**Figure 2.** Examples of each category of large language model (LLM) hallucinations along with an explanation
Note: Text highlighted in Red represents the hallucinatory generated output, while text highlighted in Blue represents user prompts or provided context that contradicts the hallucinations generated by the language model [14].

The other category is faithfulness hallucination, which deals with how well the AI's generated content aligns with what the user has asked for. Users often experience inconsistencies in three ways: instruction, context, and logical inconsistencies. Instructional inconsistency occurs when the AI's responses diverge from what the user intended, highlighting instances in which the AI misinterprets benign instructions, even if some of those deviations are related to safety measures. Context inconsistency refers to the AI contradicting the context provided by the user. Finally, logical inconsistency can arise when AI output contains contradictions, which is particularly noticeable during tasks that require reasoning.

## 3.3 Hallucination causes

Hallucinations in LLMs stem from multiple factors associated with their comprehensive capability development. The primary origins of these hallucinations are data, training, and inference. Pre-training data is crucial, as it serves as the basis for these models, enabling them to develop general skills and factual knowledge [36]. However, pre-training data can inadvertently contribute to hallucinations in LLMs through two primary mechanisms: risks associated with imperfect or unreliable data sources and the inadequate application of factual knowledge contained within the data. LLMs undergo two key phases of training: the initial pre-training phase, during which they develop general representations and acquire broad world knowledge, and a subsequent alignment phase aimed at refining their responses to align better with user instructions and preferences. Shortcomings in either of these phases can result in hallucinations in the model's outputs [14]. Decoding plays a vital role in demonstrating the capabilities of LLMs after pre-training and alignment. Therefore, faulty decoding representation and the inherent randomness of decoding strategies can lead to hallucinations in LLMs [14].

## 3.4 Hallucination mitigation

The widespread use of LLMs presents a challenge when it comes to mitigating hallucinations. Studies have proposed various mitigation methods, but some can worsen hallucinations [37, 38]. A study [37] introduced a practical approach to reducing GPT3.5 hallucinations by rectifying false information in generated sentences using retrieved knowledge. Different categories of mitigation methods were suggested, classified into the following groups:

- Fine-tuning: A Machine Learning (ML) technique is employed to fine-tune a pre-trained model for a specific context using relatively small training data [39]. A study [40] demonstrated the effectiveness of fine-tuning in mitigating LLM hallucinations, although the computational cost can be prohibitive due to LLMs' large parameter count.
- Knowledge Graphs: These methods enable the integration of structured and unstructured knowledge, providing LLMs with a broader platform to perform tasks [41]. However, designing and maintaining a curated knowledge base can require considerable time and effort.
- Memory Augmentation: Given the increasing need for deep learning techniques to expand their functionalities by incorporating additional knowledge, a recent study [42] developed an enhanced transformer specifically designed for knowledge-intensive NLP tasks. Despite the

proven benefits of memory augmentation in NLP models, no similar tests have been conducted for LLMs.
- Context Prompts: A self-monitoring prompting framework was introduced that leverages formal methods for error identification and response alignment [43]. Another study [44] introduced Self-Familiarity, a zero-resource pre-detection method designed to reduce the risk of LLMs generating inaccurate output.
- Pre-emptive Strategies: The study by Feldman et al. [45] developed a method utilizing context-tagged prompts to improve LLM response accuracy by guiding them with context prompts and validated questions.

Each approach offers unique potential for mitigating LLM hallucinations, contributing to ongoing efforts to enhance the reliability and trustworthiness of LLM-based systems.

## 4. RELATED WORK

Previous works are reviewed and organized into two sections. The First investigates hallucinations commonly encountered in various downstream tasks and explores the most prevalent hallucination types observed across these tasks. The second presents previous research exploring hallucinations in specific languages. A visual summary of this section is shown in Figure 3.
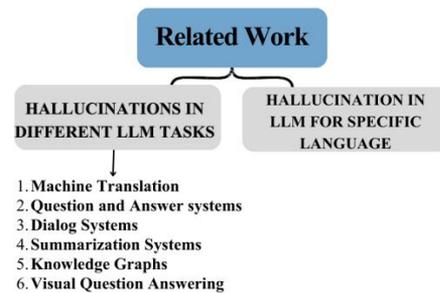


**Figure 3.** An overview of the related work section

## 4.1 Hallucinations in different Large Language Models tasks

Researchers conducted a comprehensive study investigating various types of hallucinations that may arise from different tasks performed by LLMs. The following subsections explore six specific tasks. Table 1 provides an overview of different types of hallucinations, categorized according to the respective tasks associated with LLMs.

### 4.1.1 Machine translation
Conventional translation methods verify the instances entered into the model after perturbation, as text perturbations can lead to reliable hallucinations [46, 47]. Hallucinations produced by LLMs mainly manifest as off-target translations or failed translations [23]. In situations with limited language resources, models trained on small, annotated datasets present weak performance [48]. The proliferation of pre-trained language models influences the reliability of machine translation in multilingual settings [49]. Consequently, LLMs tend to be more prone to producing hallucinations when trained on data from a single language across different scales.

### 4.1.2 Question and answer (Q&A)
Hallucination in a Q&A task occurs when an LLM provides

an incorrect answer to a user's inquiry because the LLM incorrectly infers from its source information. This can transpire even when relevant source materials are accessible. For instance, if a user inquires, "Where can a sea urchin attack a swimmer: the Andaman Sea or the Mediterranean Sea?" and the context indicates that the first option is correct, the LLM might still incorrectly respond with "Mediterranean Sea" because of its previous knowledge regarding the Mediterranean being a popular vacation spot. Instead of accurately retrieving the existing source information, the LLM might disregard the evidence and make an unwarranted inference based on its established knowledge [50]. Inaccurate external knowledge is key to producing incorrect answers, as discussed by Zheng et al. [51]. LLMs often generate incomplete but seemingly plausible responses rather than opting to provide no answer when they encounter inadequate or irrelevant information [52]. Additionally, if a system depends entirely on memorized information without referencing accessible, reliable, and accurate sources, it can result in the emergence of various types of hallucinations. A study [11] evaluates ChatGPT-3.5, ChatGPT-4.0, and Google Bard in addressing medical queries, highlighting ChatGPT-4.0's potential for accurate and comprehensive responses, whereas Google Bard achieved the lowest accuracy and showed poor responses. Moreover, scaling up models is less effective in enhancing accuracy than refining through fine-tuning with training objectives that extend beyond merely mimicking text from the web [53]. A study [53] proposed a benchmark that assesses language models' truthfulness in responding to queries. Eight hundred seventeen questions in a variety of categories were used in the study. The study tested different models like GPT-3 and T5-based models. It was discovered that, in contrast to human performance, which was accurate in 94% of cases, the best model was only true in 58% of cases.

**Table 1.** Reviewed literature of hallucination in different large language model (LLM) tasks

| Ref. | LLM Task | Dataset | Model Architecture | Hallucination Classification | Methodology |
|---|---|---|---|---|---|
| [23] | Machine Translation | WMT2018 | Encoder-Decoder | Oscillatory, largely fluent. | Natural scenario |
| [46] | Machine Translation | IWSLT2014 | Encoder-Decoder | Under perturbation, hallucination. | Source perturbation |
| [48] | Machine Translation | Wikipedia, Jig-saw, FLORES-200 | Encoder-Decoder | Full, Partial, Word-level. | Introduce pathology detection |
| [11] | Question and Answer | USMILE, Headqa, MEDMCQA, Pubmed, Medqa | Only-Decoder | Reasoning, Memory-based. | Medical Domain Hallucination Test: a custom medical benchmark |
| [51] | Question and Answer | HotpotQA, BoolQ | Only-Decoder Comprehension | Comprehension, Specificity, Factualness, Inference. | Response analysis manually |
| [52] | Question and Answer | TopiOCQA,NQ, HotpotQA | Encoder-Decoder, Decoder-only | Semantic and Symbolic analogy, Granularity, Intrinsic ambiguity discrepancies, Enumeration, Incomplete, Satisfactory Subset | Assess retrieval augmented Question and Answer. |
| [53] | Question and Answer | TruthfulQA | Encoder-Decoder, Decoder-only | mitative falsehoods | Cause imitative falsehoods |
| [54] | Dialog System | WoW | Encoder-Decoder, Decoder-only | Intrinsic, Extrinsic | Sample dialogue responses |
| [55] | Dialog System | WoW | Encoder-Decoder, Decoder-only | Generic, Uncooperativeness. | benchmark FaithDial |
| [56] | Dialog System | WoW, TopicalChat, CMU-DOG | Encoder-Decoder, Decoder-only, Encoder-only | Generic, Partial, Uncooperative. | Inferences based on the knowledge snippet |
| [57] | Summarization System | XSum, CNN/DM | Encoder-Decoder, Decoder-only | Factually inconsistent | Produce summaries from specified models |
| [58] | Summarization System | NHNet | Encoder-Decoder, Encoder-only | News headline | Majority vote of journalism specialist |
| [60] | Summarization System | MENT | Encoder-Decoder, Encoder-only | Factual, Non-factual, Intrinsic. | Identify factual objects from summaries |
| [61] | Knowledge graph | TekGen, WebNLG | Decoder-only | Subject, Object and relation hallucination | ext2KGBench: An ontology driven knowledge graph generation benchmark |
| [62] | Knowledge graph | Encyclopedic, ETC | Encoder-Decoder, Decoder-only | Knowledge | Evaluate knowledge generating ability by given facts |
| [63] | Visual Question Answer | MSCOCO | Encoder-Decoder | Object hallucination | Evaluation of caption hallucination |

### 4.1.3 Dialog system

Numerous studies have examined dialogue models as basic imitators, suggesting they merely alter data perspectives and communication patterns without generating genuinely reliable output. However, using standard benchmarks in dialog systems can actually cause models to intensify hallucinations, as shown in the study by Sun et al. [54]. In another study [55], human feedback was used to identify various types of hallucinations in Knowledge Graph (KG) grounded chatbots. Similarly, multiple experiments [54-56] conducted a meta-evaluation of hallucinations in dialogue systems that are based on knowledge by utilizing the WoW dataset.

### 4.1.4 Summarization system

LLMs enable the automatic generation of fluent abstracts, but this often comes at the cost of accuracy and faithfulness to the original content. Summaries produced by LLMs can generally be evaluated based on two types of hallucinations: intrinsic, which alter or misrepresent the information from the source, and extrinsic, which introduce new information not found in the source document [57]. Recently, extrinsic hallucinations have been getting greater focus, especially in summarization, due to LLMs' tendency to continue inputs with seemingly factual information [57, 58]. Extrinsic hallucinations are more harmful than intrinsic hallucinations because they cannot be verified from the source prompt [59]. A further refinement presented in the study by Cao et al. [60] breaks down the extrinsic category into two types: factual and non-factual. Where the former incorporates outside-world knowledge that might actually aid in understanding the text, even though it wasn't part of the original document.

### 4.1.5 Knowledge graph with Large Language Models

Text generation using a knowledge graph faces difficulties with intrinsic hallucinations, which occur when redundant or inaccurate information is introduced due to the model's reliance on its memorized knowledge [61]. A recent study [62] addressed this problem by separating correctly generated knowledge from hallucinated content. In the study by Yu et al. [62], a large LLM is proposed within a neural-symbolic framework to generate interpretable fact-checks. Hallucinations are categorized into three types: subject, relation, and object hallucinations, depending on how closely they align with the original source content.

### 4.1.6 Cross-modal system

Cross-modal tasks have made notable advancements, largely due to the exceptional language capabilities of LLMs [62]. An example of this advancement is Large Visual Language Models (LVLMs), which combine vision and natural language processing. When the original language encoder is replaced, LVLMs may persistently produce reports of items that are absent from the images, a phenomenon known as object hallucination [63]. Typically, most of these failure cases are observed in Visual QA [63] and visual Captioning tasks [64, 65].

### 4.2 Hallucination in Large Language Models for specific language

A recent study [66] introduced a benchmark called HalluQA to assess hallucinations in Chinese LLMs. This benchmark comprised 450 adversarial questions spanning different domains and categorized imitative falsehoods and factual errors as types of hallucinations. The study evaluated 24 LLMs, revealing that 18 models achieved non-hallucination rates below 50%, underscoring the considerable challenge posed by HalluQA. In a recent study [24], the authors examined hallucination tendencies in LLMs, specifically in the Bulgarian language context. They evaluated two models, namely text-DaVinci-003 and gpt-3.5-turbo-0613, utilizing the EXAMS dataset from the bgGLUE benchmark [67-69], which comprises multiple-choice questions from various subjects, including Biology, Philosophy, Geography, History, Physics, and Chemistry. The study revealed that gpt-3.5-turbo-0613 model demonstrated a higher propensity for Hallucination across all metrics assessed. Mainly, Philosophy questions obtained the highest evaluation scores due to their expansive nature, eliciting a more comprehensive range of responses compared to other subjects. Remarkably, the study identified a recurrent type of Hallucination labelled "foreign language hallucination," characterized by language-specific errors such as spelling mistakes, incorrect word order, and misused terms. These errors gained coherence when translated using machine translation tools like Google Translate. For instance, the word (rezonirane, "resonance") was interpreted as "reasoning" despite its inaccurate usage in Bulgarian.

Berbatova and Salambashev [24] encourage the widespread adoption of the latest technology in different languages. Focusing mainly on lower-resource languages can help us utilize the findings for inspiration and idea generation. Arabic language, with its unique linguistic features, is a particularly intriguing topic for research in NLP systems. As a Semitic language, Arabic boasts a complex morphological structure, utilizing a root-and-pattern system that allows the formation of numerous words from a limited set of base roots. The language is heavily inflected, with word forms varying by tense, gender, case, and number [70-72]. This complexity poses important challenges for linguistic modeling and text generation.

Moreover, the significant variation in Arabic dialects, with Modern Standard Arabic (MSA) as a formal written standard and numerous spoken dialects differ considerably in syntax, vocabulary, and phonetics [73]. This linguistic diversity, coupled with the varying data available for different dialects and the lack of an annotated data set, presents challenges in training effective language models. Considering these characteristics, research into Hallucination in Arabic language models (LLMs) is essential for improving model performance and developing effective NLP strategies, particularly for other morphologically complex and low-resource languages. Hence, this research primarily examines Hallucination in Arabic to contribute to the literature and improve the performance of Arabic LLMs.

## 5. RESEARCH METHODOLOGY

The subsequent sections present the research methodology for investigating Hallucination in Arabic LLMs for question-answering (QA) tasks with two experiments. The first one evaluated five recent and prominent LLMs, ChatGPT-4.0 and ChatGPT-3.5 [74], Gemini [75], Claude 3.5 [76], and Jais [77] using Hallucination triggering questions from the HaluEval benchmark [78]. To evaluate these models' tendency to produce hallucinations, an online prompt is utilized to pose the question in both Arabic and English. Additionally, precise criteria are specified to determine whether the response displays hallucinations. The overview of the first experiment
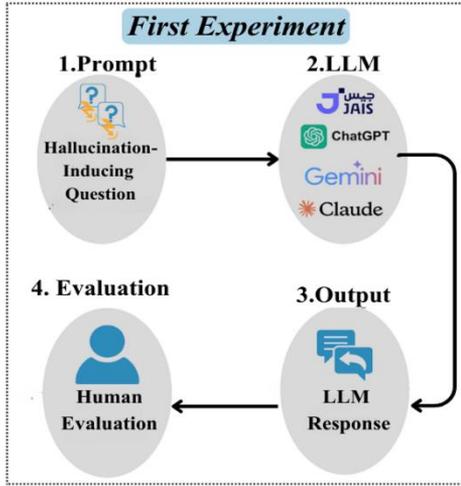
is presented in Figure 4.



**Figure 4.** First experiment's framework for evaluating hallucination in Arabic

For the second experiment, an Arabic open-domain question answering dataset is used to assess GPT-4o-mini-2024-07-18 [79]. The model selection is due to the version being a lighter variant of GPT-4o that performs efficiently while maintaining low resource consumption. Then, the selfcheckGPT framework is utilized to detect Hallucinated answers. Lastly, evaluation is performed using the BertScore metric employing the bert-base-multilingual-cased model variant to measure the similarity between the correct and generated answers. The suggested framework is shown in Figure 5. The subsequent subsection presents in detail the implementation of these two experiments.
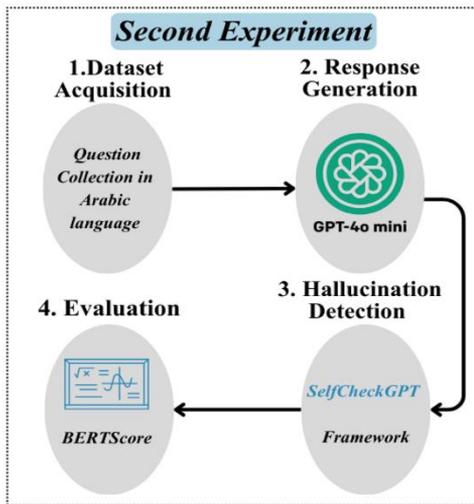


**Figure 5.** Second experiment's proposed framework for evaluating hallucination in Arabic

## 5.1 First experiment

### 5.1.1 Dataset

A set of 12 questions is meticulously selected from the HaluEval benchmark dataset [78], a comprehensive collection containing 5,000 common user questions and 30,000 domain-specific samples covering areas like QA, dialogue, and summarization that induced Hallucination in LLMs. Experiments conducted with the HaluEval benchmark shows

that not all prompts lead to hallucinated responses. A study tested and filtered the HaluEval dataset, carefully selecting 20 questions that triggered hallucinations in one or more of the chosen LLM [25]. These 20 questions were categorized into three classes: Chinese context-dependent (5 questions), English context-dependent (3 questions), and language-independent contextual (12 questions). Based on these classifications, this study extracted 12 questions from a total of 20, focusing exclusively on those categorized as language context-independent. Then, following the pattern established by the previous study to explore specific language contexts, two questions were generated that are Arabic context-dependent. Thus, questions 1 through 12 are language context-independent, while questions 13 and 14 are Arabic context-dependent. A total of 14 questions are presented in Figure 6. These questions refer to general knowledge and span multiple fields, including common sense, technology, mathematics, language analysis, geography, and science. This experiment explores the impact of language on LLMs' hallucinations by examining both the input prompt's language and the cultural context it implies. The experiment is conducted in English and Arabic, categorizing cultural background into Arabic context-dependent and language context-independent. Several factors inform this choice. First, the selected LLMs in this experiment can be categorized into Arabic LLMs, such as Jais, and multilingual LLMs. Since English is considered a high-resource language while Arabic is classified as low-resource [80, 81], this allows for a comparison of hallucination evaluation comparison based on corpus size [82]. Second, the way language is used and expressed in Arabic differs from that in English [83], which is vital to ensure the wide usability of the LLMs.



**Figure 6.** The sample of 14 questions used to induce hallucination in large language model (LLM) in the first experiments is shown in English and Arabic

### 5.1.2 Large Language Models selection

In the first experiment, five pre-trained LLMs were utilized for evaluation purposes. The LLMs used were ChatGPT-4.0 and ChatGPT-3.5 [74], developed by OpenAI; Gemini [75], developed by Google; Claude 3.5 [76], developed by Anthropic; and Jais [77], launched in 2023. Jais was specifically trained with 13 billion parameters to cater to the intricacies of the Arabic language, enhancing its understanding and generation capabilities in that linguistic context. All LLMs' calls were tested using their standard online settings, without any adjustments for specific tasks. To keep things consistent, the experiment followed a standard mid-range setup across the different models. Hence, setting the temperature value to be 1.0, the top-p value to be 0.95, and the top-k value to be 40 when applicable. The maximum number of tokens generated was determined automatically by the platform.

### 5.1.3 Evaluation metrics

The evaluation process was carefully structured to assess the tendency to hallucinate across different LLMs in the Q&A task. Each model was provided with a set of 14 questions in Arabic. The same set of questions was also presented in English for comparison. The responses were analyzed to distinguish between accurate answers and those that exhibit Hallucination. Responses were categorized as either exhibiting hallucinations or not, based on defined criteria. Given the correct answer or possible answers in case of open questions, the criteria used to determine whether an answer Exhibits Hallucination or Does Not Exhibit Hallucination are as follows:

**A. Exhibits Hallucination:** A response is classified as exhibiting Hallucination if it meets any of the following conditions:

- **Lack of Fluency**: The response is incoherent, ungrammatical, or contains excessive nonsensical content.
- **Irrelevant Correctness**: While the response contains largely accurate information, it does not directly address the question posed.
- **Inference Violation**: The response includes statements that cannot be logically inferred from the correct answer examples or contains information that contradicts them.

**B. Does Not Exhibit Hallucination:** A response is considered free from Hallucination if it satisfies either of the following conditions:

- **Logical Consistency:** The response is fully supported or can be reasonably inferred from at least one correct answer example.
- **Adherence to Undefined Cases:** If correct answer examples indicate that a question cannot be answered, then a generated response such as "I do not know" is classified as non-hallucinatory.

Since the criteria for Hallucination and non-hallucination are inherently contradictory, a response that meets the conditions for one category cannot logically belong to the other. We assess hallucination using a scoring system, where each response exhibiting Hallucination adds 1 point to the model's hallucination hallucination score. This score ranges from 0 to 14, with higher values indicating a greater tendency to generate hallucinated content. To maintain fairness, we present each question 3 times, each with a new prompt, to prevent any influence from previous dialogues. This structured

evaluation framework ensures a rigorous and objective assessment of response validity, prioritizing accuracy, relevance, and logical coherence [66].

### 5.2 Second experiment

#### 5.2.1 Dataset

The open-domain question for this experiment is utilized from the publicly available wide scale ArabicaQA dataset for open domain QA and machine reading comprehension in Arabic [84]. The authors followed carefully crafted methodology steps to create high-quality question and answer pairs, which include article selection from wiki data, question generation, filtering, and human evaluation. The data set covered many aspects, including people, dates, organizations, products, and locations. We focus on open-domain tasks, which contain more than 80,000 pairs of questions and answers divided into training, development, and testing sets. Due to resource constraints, we conducted our evaluation on a subset of 400 questions from the ArabiaQA dataset. This selection aligns with previous research that utilizes LLMs for evaluation purposes of a given task [67, 68]. The subset was selected using random sampling with a fixed random seed.

#### 5.2.2 Sampling based approach

One of the constraints of LLMs' hallucination detection approach for low-resource languages is the scarcity of annotated data with an appropriate size and format for a specific task. Furthermore, a previously suggested method requires the capability to view token probability distributions, which could be unavailable for black box models. A sampling-based approach with a reference-free framework is suggested in the previous study [85].

#### 5.2.3 Implementation details

We implement the SelfCheckGPT framework utilizing GPT-4o-mini [79]. Starting by randomly selecting 400 questions from the ArabicaQA dataset. Following that, we assess the model in a zero-evaluation manner by prompting it to generate a primary response with a temperature equal to 0 and three sample responses with a temperature equal to 1, both having a token limit of 100, resulting in a total of 1600 API calls. Bert-base-multilingual-cased variant [86] has been used to measure BertScore from the generated responses. For sentence splitting, we used the NLT sentence tokenization function instead of Spacy, which was used in the SelfCheckGPT repository to provide reliable sentence segmentation for Arabic languages.

## 6. RESULTS

### 6.1 First experiment

Five LLMs were evaluated for each of the 14 questions. The evaluation is done by one annotator to assess whether a response is hallucinated based on the evaluation criteria presented in 5.1.3. Table 2 shows the hallucination score of hallucinated responses. ChatGPT-4.0 showed superior resistance to all inducing hallucination questions in Arabic compared to its former version; ChatGPT-3.5 showed more hallucinated responses. While Gemini and Claude 3.5 responded robustly to most of the questions. Finally, Jais had the lowest resistibility among all five LLMs. Even though it

was extensively trained in the Arabic language, it provided more hallucinated responses.

**Table 2.** Comparison of hallucination hallucination scores: The number of hallucinated responses out of 14 questions for each large language model (LLM) in Arabic (AR) and English (EN)

| LLM | AR- Tendency Score | EN- Tendency Score |
|---|---|---|
| ChatGPT-4.0 | 2 | 2 |
| ChatGPT-3.5 | 5 | 4 |
| Claude 3.5 | 1 | 1 |
| Gemini | 1 | 1 |
| Jais | 8 | 7 |

According to Table 3, in all five LLMs, Hallucinated responses were triggered by questions 5, 9, 11, 13, and 14, as shown in Figure 6.

All 14 questions were asked in Arabic and English to observe the resistibility differences in the five LLMs. ChatGPT-4.0 showed the same hallucinated responses for questions 5 and 13 in both languages.

**Table 3.** The question numbers that triggered hallucinated responses out of 14 for each large language model (LLM) in Arabic and English

| LLM Model | Arabic Questions Trigger Hallucination | English Questions Trigger Hallucination |
|---|---|---|
| ChatGPT-4.0 | 5, 13 | 5, 13 |
| ChatGPT-3.5 | 4, 5, 9, 13, 14 | 5, 9, 13, 14 |
| Gemini | 13 | 14 |
| Claude 3.5 | 9 | 13 |
| Jais | 3, 4, 5, 9, 11, 12, 13, 14 | 3, 4, 7, 9, 12, 13, 14 |

In contrast, the remaining models exhibit a hallucination tendency to different question sets in each language, as listed in Table 2. Figure 7 shows Gemini's hallucinated responses in Arabic and English to question 14, which asks about the letter frequency of a given word. The question in English is, "How many times does the letter "r" occur in the word "machine"?" Gemini responded, "The word machine has two occurrences of the letter r," representing a hallucinated response. When the same question was posed in Arabic with a term that lacked the chosen letter, the response was "no occurrence. "not exhibiting Hallucination.

Another example of a hallucinated response was generated by Jais when asked about the lowest common multiple of the numbers 36 and 87. The correct answer is 1044. Figure 8 shows Jais' responses were incorrect in both languages, providing different results for each.

Figure 9 displays a snapshot of the hallucinated response generated by ChatGPT-3.5 for question 9, presented in both Arabic and English. The question asked about the probability of rolling a total of 13 with two dice, which is impossible since the numbers on the dice range from 1 to 6. The maximum total that can be achieved from rolling two dice is 12. Nevertheless, ChatGPT-3.5's response, which was consistent in both Arabic and English, incorrectly stated that rolling 6 and 7 results in a total of 13, assigning a probability of 1/36 to this event. This response is clearly a hallucination. Claude 3.5 gave the same response in both Arabic and English, as shown in Figure 10.

The most frequently hallucinated responses in both languages are triggered by question number 13, which inquires whether the names Al-Mutanabbi and Ahmed Ibn Hussein Ibn Al-Hassan Al-Ju'fi Al-Kindi Al-Kufi refer to the same individual. The correct answer is that they do, yet most models cannot accurately identify the relationship between these two names, even though the model can accurately identify the full name of Al-Mutanabbi when asked.

**6.2 Second experiment**

In this section, we will go through the experimental findings to illustrate the performance of the suggested technique and the hallucination result. We examined the ability of the GPT-4o-mini model on open-domain questions for Arabic languages. The final hallucination results are calculated by comparing each sentence in the main response with each sampled response. The average hallucination score for the GPT-4o-mini is 19.63%. The lowest score is equal to 0, indicating that hallucinations are not observed, while the highest generated score is 33%.

Figure 11 shows a sample of responses. Question number one achieved the highest hallucination score among all questions. The main response generated a factual response, but others provided vague or Hallucination responses.

The first and third sample responses explicitly state that they lack the required knowledge to provide an answer, making them less informative while effectively avoiding hallucinated responses. Sample two responded with hallucination information contradicting the main answer, leading to factual inconsistency.

For the second question, the response indicates that the model can identify the general information about the place of birth: "Egypt". However, it fails to provide more detailed information about the city, resulting in a unique city for each response. Furthermore, sample 2 provided factual information unrelated to the question. This response achieved an 18.83% hallucination score. In contrast, question three yielded a response of a zero hallucination score, demonstrating consistency across all answers.
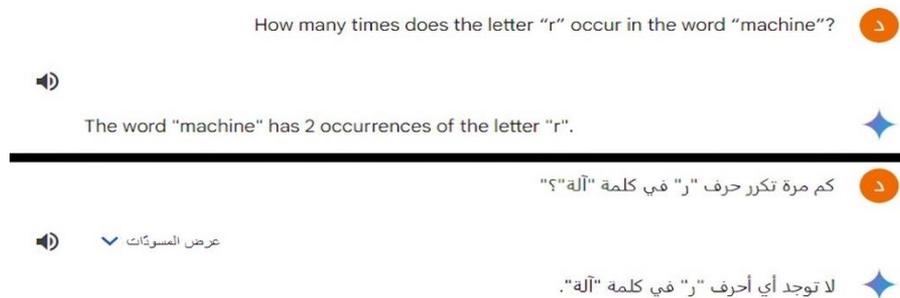


How many times does the letter "r" occur in the word "machine"?

The word "machine" has 2 occurrences of the letter "r".

كم مرة تكرر حرف "ر" في كلمة "آلة"؟

عرض المسوّدات ⌄

لا توجد أي أحرف "ر" في كلمة "آلة".

**Figure 7.** A screenshot of the generated response from Gemini to question 14 in Arabic and English

**Figure 8.** A screenshot of the generated response from Jais to question 4 in Arabic and English



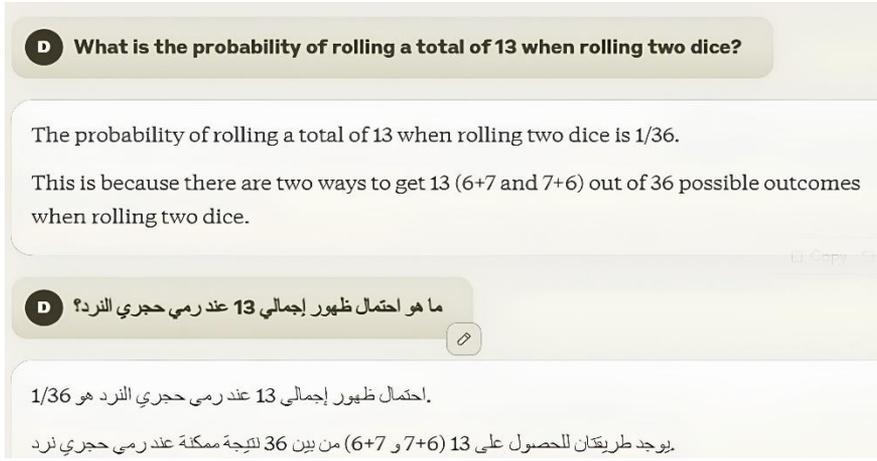**Figure 9.** A screenshot of the generated response from ChatGPT-3.5 to question 9 in Arabic and English

46

**Figure 10.** A screenshot of the generated response from Claude 3.5 to question 9 in Arabic and English

| | Question | Response Type | Answer |
|---|---|---|---|
| 1 | ما هي الدرجات العلمية التي حصل عليها مسعد العودي؟ | Main | محمد مسعد العودي هو أكاديمي وباحث معروف في مجاله. حصل على عدة درجات علمية، منها: 1. **البكالوريوس** في اللغة العربية. 2. **الماجستير** في الأدب العربي.3. **الدكتوراه** في الأدب العربي، إذا كنت بحاجة إلى معلومات أكثر تفصيلاً حول مسيرته الأكاديمية أو إنجازاته، فلا تتردد في طرح المزيد من الأسئلة! |
| | | Sample 1 | لا أملك معلومات دقيقة حول محمد مسعد العودي بشكل شخصي، بما في ذلك تفاصيل عن درجاته العلمية. قد تحتاج إلى البحث في مصادر أخرى أو المواقع الرسمية للحصول على معلومات دقيقة عن سيرته الذاتية والأكاديمية. إذا كان لديك أي استفسارات أخرى أو تحتاج معلومات عن موضوع آخر، فلا تتردد في طرح ذلك! |
| | | Sample 2 | محمد مسعد العودي هو أكاديمي وباحث سعودي، حصل على درجاته العلمية في مجالات متعددة. حصل على درجة البكالوريوس في الكيمياء، ودرجة الماجستير في الكيمياء كذلك، من جامعة الملك سعود. ثم حصل على درجة الدكتوراه في مجال الكيمياء من جامعة تكساس A&M في الولايات المتحدة. |
| | | Sample 3 | محمد مسعد العودي هو شخصية أكاديمية بارزة، وقد حصل على مجموعة من الدرجات العلمية. إذا كنت تبحث عن معلومات دقيقة حول درجاته العلمية، أو ضرورياً أن نعرف انتماءاته الأكاديمية، يفضل الرجوع إلى مصادر موثوقة مثل المواقع الأكاديمية الرسمية أو السيرة الذاتية المحدثة له. مؤخراً قد لا تتوفر لي معلومات دقيقة بعد أكتوبر 2021، لذا قد تحتاج إلى التحقق من مصادر أخرى للحصول على أحدث المعلومات. |
| 2 | أين ولد عبد الوهاب محمد عبد الحافظ؟ | Main | عبد الوهاب محمد عبد الحافظ وُلد في محافظة الدقهلية في مصر. إذا كنت بحاجة إلى مزيد من المعلومات حوله أو حول إنجازاته، فلا تتردد في طرح سؤالك! |
| | | Sample 1 | عبد الوهاب محمد عبد الحافظ ولد في محافظة المنيا بمصر. إذا كنت بحاجة إلى معلومات إضافية عنه، فلا تتردد في السؤال! |
| | | Sample 2 | عبد الوهاب محمد عبد الحافظ هو شخصية أكاديمية وباحث مصري معروف في مجالات العلوم الإنسانية والاجتماعية، ولكنه ليس شخصية مشهورة عالمياً مثل الفنانين أو الرياضيين. إذا كان لديك سياق محدد حول الشخص الذي تسأل عنه، يرجى توضيحه لتتمكن من مساعدتك بشكل أفضل. |
| | | Sample 3 | عبد الوهاب محمد عبد الحافظ وُلد في مدينة كفر الشيخ في مصر. إذا كان لديك أي أسئلة أخرى أو تحتاج إلى معلومات إضافية، فلا تتردد في طرحها! |
| 3 | ما هي عاصمة أستراليا؟ | Main | عاصمة أستراليا هي كانبيرا. |
| | | Sample 1 | عاصمة أستراليا هي كانبيرا. |
| | | Sample 2 | عاصمة أستراليا هي كانبيرا. |
| | | Sample 3 | عاصمة أستراليا هي كانبيرا. |

**Figure 11.** Example of generated response that contains the main answer and three sample answers for a given question
Note: The Non-factual statements are highlighted in red. English translations are shown in Figure A1.

## 6.3 Comparative analysis

Comparing the two experiments provides a clearer picture of how hallucinations occur in LLMs. The first experiment examined differences across models and languages using carefully designed questions intended to trigger hallucinations. In contrast, the second experiment explored how consistent a single model's responses are when generating multiple outputs for the same question. The findings show clear variation in hallucination behavior depending on both the model and the language used. ChatGPT-4.0 demonstrated greater resistance to hallucination than ChatGPT-3.5, reflecting improvements in newer model versions. Gemini and Claude 3.5 generally performed well, although their responses were influenced by language. Jais, on the other hand, showed the weakest performance, producing hallucinated answers in both Arabic and English.

The second experiment, which focused on GPT-4o-mini, revealed an average hallucination rate of 19.63%. While some questions produced fully consistent and accurate responses, others resulted in noticeable inconsistencies. This variation indicates that hallucinations can depend strongly on the type of question, even within the same model. Overall, the combined results suggest that hallucinations are shaped by both model design and prompt characteristics. Although advances in model architecture help reduce hallucinations,

they do not fully eliminate the problem, highlighting the importance of evaluating models both across different systems and within individual models to better assess their reliability.

## 6.4 Risk scenario

Hallucinations in LLMs present significant risks when their outputs are used for decision-making, educational purposes, or knowledge dissemination [5, 12, 31]. Furthermore, beyond domain-specific risks, hallucination may also present societal implications that may lead to reduced public trust in AI technologies [87]. In the context of this study, users of Arabic and English question-answer systems, including students, educators, and professionals, may inadvertently receive inaccurate or misleading information. Factual errors, such as incorrect calculations, misrepresented historical details, or inaccurate scientific information, can lead to misunderstandings and the unintentional spread of misinformation. Misinterpretation of names, entities, or culturally specific references may produce outputs that conflict with the intended context, potentially undermining user trust in the model. Additionally, inconsistencies in following instructions, such as generating overly long responses or misaligned sentence structures, can reduce the clarity and usability of model outputs. These risks are particularly pronounced in underrepresented languages like Arabic, where limited training data may increase the likelihood of hallucinations. The potential consequences range from minor confusion in educational contexts to more serious issues when models are used for professional or public decision support. Recognizing these risks highlights the importance of both cross-model and intra-model evaluations, as well as the implementation of validation mechanisms, to mitigate the impact of hallucinations in real-world applications.

## 6.5 Discussion

This research contributes by addressing the gap and investigating hallucinations in LLMs in the Arabic language. Two experiments were conducted focusing on the question-answer task. The findings of this research are significant, as they shed light on the tendency to produce hallucinated responses of different LLMs in Arabic and English. Claude 3.5 and Gemini outperformed all other LLMs in the first experiment, successfully resisting 13 out of 14 questions. ChatGPT-4.0 ranked second, generating hallucinated responses for only two questions. ChatGPT-4.0 displayed consistent performance across Arabic and English when responding to the same 14 questions, with no significant differences observed between the two languages. Most hallucinated responses were related to mathematics questions (such as questions 4, 7, and 9) and Arabic-contextual questions (like questions 13 and 14).

Based on hallucination categories in the related work section and our observations, hallucinated responses to question 3 can be classified as instruction inconsistency. In this case, the model was asked to generate five-word sentences to describe a person, but produced five longer sentences instead. Question 4 involved finding the lowest common multiple of two numbers, and the hallucinated response here fell under logical inconsistency. Similarly, hallucinated responses to questions 7, 9, and 14 also belong to this category. For example, when an LLM claims that a six-sided die shows

a 7, this could be considered a context inconsistency and intrinsic Hallucination, as the model generates data that contradicts the provided context.

For question 5, it is essential to note that 'G' is not a valid digit in the hexadecimal number system; therefore, there is no associated color. However, the model generated a response related to the RGB color system, which can be interpreted as a factual fabrication. A similar issue arises with question 11, where the model inaccurately suggests that red capes irritate bulls; in fact, the use of this color has historical significance. The model fabricated unrelated facts to create a justification for the question. Finally, question 13 highlighted issues with understanding the Arabic naming system, which many LLMs struggled to grasp.

This result highlights the importance of language structure and how models process each language differently due to the size of the training data or inference process. These findings underscore the urgent need for language models to be more prepared to handle the intricacies of languages other than English. Literature has demonstrated that LLMs are more likely to produce hallucinatory responses in non-English languages such as Chinese [66] and Japanese [88]. Furthermore, the hallucinated responses in English, despite the extensive training data, suggest that training data is not enough to mitigate hallucinations; improvements in inference processes are also necessary.

For the second experiment, the SelfcheckGPT approach with the BERTScore variant has been used to investigate hallucination tendencies in the GPT-4o-mini model for open Arabic questions. From the results, it was noted that when the model struggles to generate accurate answers, particularly for questions about people, it tends to produce overly long responses. Additionally, in some cases, even though all generated answers were inconsistent with each other and failed to provide a precise answer, they still received low hallucination scores. This may be attributed to the presence of common phrases shared among the responses, such as repeating the same information provided in the question, which does not contribute to delivering the required answer.

## 7. CONCLUSION

This study explored hallucinations, where LLMs generate fabricated or inconsistent responses, focusing on the QA task in the Arabic language and presenting a thorough overview of the definition, classification, causes, and mitigation of LLM hallucinations. To examine the phenomenon of LLMs' hallucination in question-answering tasks within the Arabic language, this study undertook two experiments. The first experiment was predicated on zero-shot learning, employing LLM models such as ChatGPT-3.5 and Gemini. The second experiment implemented a sampling-based approach by utilizing the selfcheckGPT framework. The findings of the first experiment shed light on the tendency of five prominent LLMs to produce hallucinated responses in Arabic and English. The second experiment applied the BERT Score to measure the similarity between the model's responses, and lower scores indicate hallucinated responses. Overall, this study provides insight by investigating the level of Hallucination of LLMs for the Arabic language, which poses unique challenges compared to other languages, such as morphological richness, script variations, and lack of annotated datasets [70, 71].

These challenges are characteristic of many low-resource languages, where limited high-quality training data, scarce benchmarks, and diverse dialectal variations hinder the development of reliable language models. Morphologically complex languages, including Arabic, often require models to understand intricate root-and-pattern systems, extensive inflectional paradigms, and context-dependent word forms, which can increase the likelihood of hallucinated outputs. Additionally, the scarcity of annotated corpora for both formal and dialectal varieties limits opportunities for model fine-tuning and evaluation, further exacerbating performance gaps compared to high-resource languages like English. As observed in studies of other low-resource languages, such as Chinese [66] and Bulgarian [24], hallucination rates tend to be higher, particularly in domains that require factual precision or context-sensitive reasoning. By focusing on Arabic, this research highlights the broader need for improved resources, evaluation benchmarks, and targeted strategies for low-resource languages, contributing to more accurate and reliable NLP systems across linguistically diverse contexts.

A key limitation lies in the experiment's scale. Testing only a few LLMs on a small dataset limits the generalizability of the findings. Different languages possess unique grammatical structures, expressions, and cultural contexts that can influence hallucinations in varying ways. Future work should involve testing a wider range of LLMs across diverse languages and larger datasets to gain a more comprehensive understanding. This would allow researchers to explore the influence of language on hallucinations, delve deeper into the root causes, and develop methods for improving LLM robustness.

## DECLARATIONS

All authors declare that they have no conflicts of interest.

## REFERENCES

[1] Güler, A., Akgül, İ. (2022). A review on the science of natural language processing.

[2] Thapa, S., Adhikari, S. (2023). ChatGPT, bard, and large language models for biomedical research: Opportunities and pitfalls. Annals of Biomedical Engineering, 51(12): 2647-2651. https://doi.org/10.1007/S10439-023-03284-0

[3] Caruccio, L., Cirillo, S., Polese, G., Solimando, G., Sundaramurthy, S., Tortora, G. (2024). Can ChatGPT provide intelligent diagnoses? A comparative study between predictive models and ChatGPT to define a new medical diagnostic bot. Expert Systems with Applications, 235: 121186. https://doi.org/10.1016/j.eswa.2023.121186

[4] Gu, Q. (2023). LLM-based code generation method for Golang compiler testing. In Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, pp. 2201-2203. https://doi.org/10.1145/3611643.3617850

[5] Meyer, J.G., Urbanowicz, R.J., Martin, P.C., O'Connor, K., Li, R., Peng, P.C., Bright, T.J., Tatonetti, N., Won, K.J., Gonzalez-Hernandez, G., Moore, J.H. (2023). ChatGPT and large language models in academia: Opportunities and challenges. BioData Mining, 16(1): 20. https://doi.org/10.1186/S13040-023-00339-9

[6] Fatouros, G., Soldatos, J., Kouroumali, K., Makridis, G., Kyriazis, D. (2023). Transforming sentiment analysis in the financial domain with ChatGPT. Machine Learning with Applications, 14: 100508. https://doi.org/10.1016/j.mlwa.2023.100508

[7] Qiu, Y. (2024). The impact of LLM hallucinations on motor skill learning: A case study in badminton. IEEE Access, 12: 139669-139682. https://doi.org/10.1109/ACCESS.2024.3444783

[8] Merriam Webster. (2025). Definition of HALLUCINATION, Merriam-Webster.com. https://www.merriam-webster.com/dictionary/hallucination, accessed on 2 Feb 2025.

[9] Perković, G., Drobnjak, A., Botički, I. (2024). Hallucinations in LLMS: Understanding and addressing challenges. In 2024 47th MIPRO ICT and Electronics Convention (MIPRO), Opatija, Croatia, pp. 2084-2088. https://doi.org/10.1109/MIPRO60963.2024.10569238

[10] Reddy, G.P., Pavan Kumar, Y.V., Prakash, K.P. (2024). Hallucinations in large language models (LLMs). In Proceedings of the 2024 IEEE Open Conference of Electrical, Electronic and Information Sciences (eStream), Vilnius, Lithuania, pp. 1-6. https://doi.org/10.1109/eStream61684.2024.10542617

[11] Pal, A., Umapathi, L.K., Sankarasubbu, M. (2023). Med-HALT: Medical domain hallucination test for large language models. In Proceedings of the Conference on Computational Natural Language Learning. https://doi.org/10.48550/ARXIV.2307.15343

[12] Raj, R., Singh, A., Kumar, V., Verma, P. (2023). Analyzing the potential benefits and use cases of ChatGPT as a tool for improving the efficiency and effectiveness of business operations. BenchCouncil Transactions on Benchmarks, Standards and Evaluations, 3(3): 100140. https://doi.org/10.1016/j.tbench.2023.100140

[13] Zhou, H., Hu, C., Yuan, Y., Cui, Y., et al. (2024). Large language model (LLM) for telecommunications: A comprehensive survey on principles, key techniques, and opportunities. IEEE Communications Surveys & Tutorials, 27(3): 1955-2005.

[14] Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Qin, B., Liu, T. (2025). A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. ACM Transactions on Information Systems, 43(2): 1-55. https://doi.org/10.1145/3703155

[15] Chui, M., Hazan, E., Roberts, R., Singla, A., Smaje, K. (2023). The economic potential of generative AI: The next productivity frontier. McKinsey & Company. https://cloudeurope.nl/images/Downloads/the-economic-potential-of-generative-ai-the-next-productivity-frontier.pdf.

[16] Plevris, V., Papazafeiropoulos, G., Jiménez Rios, A. (2023). Chatbots put to the test in math and logic problems: A comparison and assessment of ChatGPT-3.5, ChatGPT-4, and Google Bard. AI, 4(4): 949-969. https://doi.org/10.3390/ai4040048

[17] Richter, F. (2022). Infographic: English Is the Internet's Universal Language. Statista Info-graphics. Statista. https://www.statista.com/chart/26884/languages-on-the-internet/.

[18] Al-Tamimi, A.K., Bani-Isaa, E., Al-Alami, A. (2021). Active learning for Arabic text classification. In 2021 International Conference on Computational Intelligence and Knowledge Economy (ICCIKE), Dubai, United Arab Emirates, pp. 123-126. https://doi.org/10.1109/ICCIKE51210.2021.9410758

[19] International Telecommunication Union (ITU) World Telecommunication/ICT Indicators Database (2022) 'Individuals using the internet (% of population) - Arab World | Data', data.worldbank.org. https://data.worldbank.org/ indicator/IT.NET.USER.ZS?locations=1A, accessed on Feb. 03, 2025.

[20] Senator, F., Lakhfif, A., Zenbout, I., Boutouta, H., Mediani, C. (2025). Leveraging ChatGPT for enhancing Arabic NLP: Application for semantic role labeling and cross-lingual annotation projection. IEEE Access, 13: 3707-3725. https://doi.org/10.1109/ACCESS.2025.3525493

[21] Chen, C., Shu, K. (2023). Can LLM-generated misinformation be detected?. arXiv:2309.13788. http://arxiv.org/abs/2309.13788

[22] Chen, Z., Jiang, F., Chen, J., Wang, T., Yu, F., Chen, G., Zhang, H., Liang, J., Zhang, C., Zhang, Z., Li, J., Wan, X., Wang, B., Li, H. (2023). Phoenix: Democratizing ChatGPT across languages. arXiv:2304.10453. https://arxiv.org/abs/2304.10453v1

[23] Guerreiro, N.M., Alves, D.M., Waldendorf, J., Haddow, B., Birch, A., Colombo, P., Martins, A.F. (2023). Hallucinations in large multilingual translation models. Transactions of the Association for Computational Linguistics, 11: 1500-1517. https://doi.org/10.1162/tacl_a_00615

[24] Berbatova, M., Salambashev, Y. (2023). Evaluating hallucinations in large language models for Bulgarian language. In Proceedings of the 8th Student Research Workshop associated with the International Conference Recent Advances in Natural Language Processing, pp. 55-63.

[25] Hu, R., Zhong, J., Ding, M., Ma, Z., Chen, M. (2023). Evaluation of hallucination and robustness for large language models. In 2023 IEEE 23rd International Conference on Software Quality, Reliability, and Security Companion (QRS-C), Chiang Mai, Thailand, pp. 374-382. https://doi.org/10.1109/QRS-C60940.2023.00089

[26] Sadat, M., Zhou, Z., Lange, L., Araki, J., Gundroo, A., Wang, B., Menon, R.R., Parvez, M.R., Feng, Z. (2023). DelucionQA: Detecting hallucinations in domain-specific question answering. arXiv:2312.05200. http://arxiv.org/abs/2312.05200

[27] Das, S., Chatterji, S., Mukherjee, I. (2023). Combating Hallucination and Misinformation: Factual Information Generation with Tokenized Generative Transformer. In Proceedings of the Joint 3rd International Conference on Natural Language Processing for Digital Humanities and 8th International Workshop on Computational Linguistics for Uralic Languages, pp. 143-152.

[28] Huang, H., Tang, T., Zhang, D., Zhao, X., Song, T., Xia, Y., Wei, F. (2023). Not all languages are created equal in LLMs: Improving multilingual capability by cross-lingual-thought prompting. In Findings of the Association for Computational Linguistics: EMNLP 2023, pp. 12365-12394.

https://doi.org/10.18653/v1/2023.findings-emnlp.826

[29] Bang, Y. (2023). A multitask, multilingual, multimodal evaluation of ChatGPT on reasoning, hallucination, and interactivity. arXiv preprint arXiv:2302.04023.

[30] Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y.J., Madotto, A., Fung, P. (2023). Survey of hallucination in natural language generation. ACM Computing Surveys, 55(12): 1-38. https://doi.org/10.1145/3571730

[31] Yilmaz, R., Yilmaz, F.G.K. (2023). Augmented intelligence in programming learning: Examining student views on the use of ChatGPT for programming learning. Computers in Human Behavior: Artificial Humans, 1(2): 100005. https://doi.org/10.1016/j.chbah.2023.100005

[32] Li, J., Dada, A., Puladi, B., Kleesiek, J., Egger, J. (2024). ChatGPT in healthcare: A taxonomy and systematic review. Computer Methods and Programs in Biomedicine, 245: 108013. https://doi.org/10.1016/j.cmpb.2024.108013

[33] Rawte, V., Chakraborty, S., Pathak, A., Sarkar, A., Tonmoy, S.T.I., Chadha, A., Sheth, A., Das, A. (2023). The troubling emergence of hallucination in large language models-an extensive definition, quantification, and prescriptive remediations. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 2541-2573. https://doi.org/10.18653/v1/2023.emnlp-main.155

[34] Chen, Y., Fu, Q., Yuan, Y., Wen, Z., Fan, G., Liu, D., Zhang, D., Li, Z., Xiao, Y. (2023). Hallucination detection: Robustly discerning reliable answers in large language models. In Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, pp. 245-255. https://doi.org/10.1145/3583780.3614905

[35] Krasadakis, P., Sakkopoulos, E., Verykios, V.S. (2024). A survey on challenges and advances in natural language processing with a focus on legal informatics and low-resource languages. Electronics, 13(3): 648. https://doi.org/10.3390/electronics13030648

[36] Zhou, Y., Cui, C., Yoon, J., Zhang, L., Deng, Z., Finn, C., Bansal, M., Yao, H. (2023). Analyzing and mitigating object hallucination in large vision-language models. arXiv preprint arXiv:2310.00754. https://doi.org/10.48550/arXiv.2310.00754

[37] Varshney, N., Yao, W., Zhang, H., Chen, J., Yu, D. (2023). A stitch in time saves nine: Detecting and mitigating hallucinations of LLMS by validating low-confidence generation. arXiv preprint arXiv:2307.03987. https://doi.org/10.48550/arXiv.2307.03987

[38] Kazlaris, I., Antoniou, E., Diamantaras, K., Bratsas, C. (2025). From illusion to insight: A taxonomic survey of hallucination mitigation techniques in LLMs. AI, 6(10): 260. https://doi.org/10.20944/preprints202508.1942.v1

[39] Church, K.W., Chen, Z., Ma, Y. (2021). Emerging trends: A gentle introduction to fine-tuning. Natural Language Engineering, 27(6): 763-778. https://doi.org/10.1017/S1351324921000322

[40] Lee, C., Cho, K., Kang, W. (2019). Mixout: Effective regularization to finetune large-scale pretrained language models. arXiv preprint arXiv:1909.11299. https://doi.org/10.48550/arXiv.1909.11299

[41] Moiseev, F., Dong, Z., Alfonseca, E., Jaggi, M. (2022) SKILL: Structured knowledge infusion for large

language models. In NAACL 2022 - 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, pp. 1581-1588. https://doi.org/10.18653/v1/2022.naacl-main.113

[42] Wu, Y., Zhao, Y., Hu, B., Minervini, P., Stenetorp, P., Riedel, S. (2022). An efficient memory-augmented transformer for knowledge-intensive NLP tasks. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, pp. 5184-5196. https://doi.org/10.18653/v1/2022.emnlp-main.346

[43] Jha, S., Jha, S.K., Lincoln, P., Bastian, N.D., Velasquez, A., Neema, S. (2023). Dehallucinating large language models using formal methods guided iterative prompting. In 2023 IEEE International Conference on Assured Autonomy (ICAA), Laurel, MD, USA, pp. 149-152. https://doi.org/10.1109/ICAA58325.2023.00029

[44] Feldman, P., Foulds, J.R., Pan, S. (2023). Zero-resource hallucination prevention for large language models. arXiv:2309.02654.
https://doi.org/10.48550/arXiv.2309.02654

[45] Feldman, P., Foulds, J.R., Pan, S. (2023). Trapping LLM hallucinations using tagged context prompts. arXiv:2306.06085.
https://doi.org/10.48550/arXiv.2306.06085

[46] Raunak, V., Menezes, A., Junczys-Dowmunt, M. (2021). The curious case of hallucinations in neural machine translation. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 1172-1183. https://doi.org/10.18653/v1/2021.naacl-main.92

[47] Bawden, R., Yvon, F. (2023). Investigating the translation performance of a large multilingual language model: The case of BLOOM. arXiv:2303.01911. https://doi.org/10.48550/ARXIV.2303.01911

[48] Dale, D., Voita, E., Lam, J., Hansanti, P., Ropers, C., Kalbassi, E., Gao, C., Barrault, L., Costa-jussà, M.R. (2023). HalOmi: A manually annotated benchmark for multilingual hallucination and omission detection in machine translation. arXiv:2305.11746. https://doi.org/10.48550/ARXIV.2305.11746

[49] Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., Stoyanov, V. (2020). Unsupervised cross-lingual representation learning at scale. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 8440-8451. https://doi.org/10.18653/v1/2020.acl-main.747

[50] Galitsky, B. (2025). Truth-o-meter: Collaborating with LLM in fighting its hallucinations. In Interdependent Human-Machine Teams, pp. 175-210. Academic Press. https://doi.org/10.1016/B978-0-443-29246-0.00004-3

[51] Zheng, L.M., Chiang, W.L., Sheng, Y., Zhuang, S.Y., et al. (2023). Judging LLM-as-a-judge with MT-bench and Chatbot Arena. Advances in Neural Information Processing Systems, 36: 46595-46623.

[52] Adlakha, V., BehnamGhader, P., Lu, X.H., Meade, N., Reddy, S. (2023). Evaluating correctness and faithfulness of instruction-following models for question answering. arXiv:2307.16877. https://doi.org/10.48550/ARXIV.2307.16877

[53] Lin, S., Hilton, J., Evans, O. (2022). Truthfulqa: Measuring how models mimic human falsehoods. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 3214-3252. https://doi.org/10.18653/v1/2022.acl-long.229

[54] Sun, W., Shi, Z., Gao, S., Ren, P., de Rijke, M., Ren, Z. (2023). Contrastive learning reduces hallucination in conversations. Proceedings of the AAAI Conference on Artificial Intelligence, 37(11): 13618-13626. https://doi.org/10.1609/aaai.v37i11.26596

[55] Dziri, N., Kamalloo, E., Milton, S., Zaïane, O.R., Yu, M., Ponti, E.M., Reddy, S. (2022). Faithdial: A faithful benchmark for information-seeking dialogue. Transactions of the Association for Computational Linguistics, 10: 1473-1490. https://doi.org/10.1162/tacl_a_00529

[56] Dziri, N., Rashkin, H., Linzen, T., Reitter, D. (2022). Evaluating attribution in dialogue systems: The BEGIN benchmark. Transactions of the Association for Computational Linguistics, 10: 1066-1083. https://doi.org/10.1162/tacl_a_00506

[57] Tam, D., Mascarenhas, A., Zhang, S., Kwan, S., Bansal, M., Raffel, C. (2023). Evaluating the Factual Consistency of Large Language Models Through News Summarization. In Proceedings of the Annual Meeting of the Association for Computational Linguistics, pp. 5220-5255. https://doi.org/10.18653/V1/2023.FINDINGS-ACL.322

[58] Shen, J., Liu, J., Finnie, D., Rahmati, N., Bendersky, M., Najork, M. (2023). "Why is this misleading?": Detecting News Headline Hallucinations with Explanations. In Proceedings of the ACM Web Conference 2023, pp. 1662-1672. https://doi.org/10.1145/3543507.3583375

[59] Rawte, V., Chakraborty, S., Pathak, A., Sarkar, A., et al. (2023). The troubling emergence of hallucination in large language models—An extensive definition, quantification, and prescriptive remediations. arXiv: 2310.04988

[60] Cao, M., Dong, Y., Cheung, J.C.K. (2022). Hallucinated but factual! inspecting the factuality of hallucinations in abstractive summarization. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 3340-3354. https://doi.org/10.18653/v1/2022.acl-long.236

[61] Mihindukulasooriya, N., Tiwari, S., Enguix, C.F., Lata, K. (2023). Text2KGBench: A benchmark for ontology-driven knowledge graph generation from text. arXiv:2308.02357.
https://doi.org/10.48550/ARXIV.2308.02357

[62] Yu, J., Wang, X., Tu, S., Cao, S., et al. (2023). KoLA: Carefully benchmarking world knowledge of large language models. arXiv:2306.09296. https://doi.org/10.48550/ARXIV.2306.09296

[63] Li, Y., Du, Y., Zhou, K., Wang, J., Zhao, W.X., Wen, J. (2023). Evaluating object hallucination in large vision-language models. arXiv:2305.10355. https://doi.org/10.48550/ARXIV.2305.10355

[64] Biten, A.F., Gomez, L., Karatzas, D. (2022). Let there be a clock on the beach: Reducing object hallucination in image captioning. In Proceedings - 2022 IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2022, Waikoloa, HI, USA, pp. 2473-2482. https://doi.org/10.1109/WACV51458.2022.00253

[65] Favero, A., Zancato, L., Trager, M., Choudhary, S.,

Perera, P., Achille, A., Swaminathan, A., Soatto, S. (2024). Multi-modal hallucination control by visual information grounding. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14303-14312.

[66] Cheng, Q., Sun, T., Zhang, W., Wang, S., Liu, X., Zhang, M., He, J., Huang, M., Yin, Z., Chen, K., Qiu, X. (2023). Evaluating hallucinations in Chinese large language models. https://doi.org/10.48550/arXiv.2310.03368

[67] Al-Thubaity, A., Alkhereyf, S., Murayshid, H., Alshalawi, N., Omirah, M., Alateeq, R., Almutairi, R., Alsuwailem, R., Alhassoun, M., Alkhanen, I. (2023). Evaluating ChatGPT and bard AI on Arabic sentiment analysis. In Proceedings of ArabicNLP 2023, pp. 335-349. https://doi.org/10.18653/v1/2023.arabicnlp-1.27

[68] Muller, S., Loison, A., Omrani, B., Viaud, G. (2024) GROUSE: A benchmark to evaluate evaluators in grounded question answering. arXiv:2409.06595. https://arxiv.org/abs/2409.06595

[69] Hardalov, M., Atanasova, P., Mihaylov, T., Angelova, G., Simov, K., Osenova, P., Stoyanov, V., Koychev, I., Nakov, P., Radev, D. (2023). bgGLUE: A Bulgarian general language understanding evaluation benchmark. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 8733-8759. https://doi.org/10.18653/v1/2023.acl-long.487

[70] AlAfnan, M.A. (2025). Artificial Intelligence and language: Bridging Arabic and English with technology. Journal of Ecohumanism, 4(1): 240-256. https://doi.org/10.62754/joe.v3i8.4961

[71] Khalatia, M.M., Al-Romanyb, T.A.H. (2020). Artificial intelligence development and challenges (Arabic language as a model). Artificial Intelligence, 13(5): 916-926.

[72] Sallam, M., Al-Mahzoum, K., Almutawaa, R.A., Alhashash, J.A., Dashti, R.A., AlSafy, D.R., Almutairi, R.A., Barakat, M. (2024). The performance of OpenAI ChatGPT-4 and Google Gemini in virology multiple-choice questions: A comparative analysis of English and Arabic responses. BMC Research Notes, 17(1): 247. https://doi.org/10.1186/s13104-024-06920-7

[73] Al-Azani, S., Alturayeif, N., Abouelresh, H., Alhunief, A. (2024). A comprehensive framework and empirical analysis for evaluating large language models in Arabic dialect identification. In 2024 International Joint Conference on Neural Networks (IJCNN), Yokohama, Japan, pp. 1-7. https://doi.org/10.1109/IJCNN60899.2024.10651099

[74] OpenAI. (2015). GPT-4, Openai.com. https://openai.com/ index/gpt-4/.

[75] Google DeepMind. (2023). Gemini - Google DeepMind, deepmind.google. https://deepmind.google/technologies/gemini/.

[76] Anthropic. (2023). 'Claude', Claude.ai. https://claude.ai/ chats.

[77] Jais AI. (2023). Advanced Arabic Language AI model. Jais Ai. https://jais.pro/.

[78] Li, J., Cheng, X., Zhao, X., Nie, J.Y., Wen, J.R. (2023). HaluEval: A large-scale hallucination evaluation benchmark for large language models. Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 6449-6464. https://doi.org/10.18653/v1/2023.emnlp-main.397

[79] OpenAI. (2024). GPT-4o mini: Advancing cost-efficient intelligence, Openai.com. https://openai.com/index/ gpt-4o-mini-advancing-cost-efficient-intelligence/.

[80] Kholodna, N., Julka, S., Khodadadi, M., Gumus, M.N., Granitzer, M. (2024). LLMs in the Loop: Leveraging Large Language Model Annotations for Active Learning in Low-Resource Languages. In: Bifet, A., Krilavičius, T., Miliou, I., Nowaczyk, S. (eds) Machine Learning and Knowledge Discovery in Databases. Applied Data Science Track. ECML PKDD 2024. Lecture Notes in Computer Science(), vol 14950. Springer, Cham. https://doi.org/10.1007/978-3-031-70381-2_25

[81] Bommasani, R. (2021). On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258.

[82] Alammary, A.S. (2025). Investigating the impact of pretraining corpora on the performance of Arabic BERT models. The Journal of Supercomputing, 81(1): 187. https://doi.org/10.1007/s11227-024-06698-2

[83] Essam, M., Deif, M., Elgohary, R. (2024). Deciphering Arabic question: A dedicated survey on Arabic question analysis methods, challenges, limitations and future pathways. Artificial Intelligence Review, 57(9): 251. https://doi.org/10.1007/s10462-024-10880-6

[84] Abdallah, A., Kasem, M., Abdalla, M., Mahmoud, M., Elkasaby, M., Elbendary, Y., Jatowt, A. (2024) 'ArabicaQA: A comprehensive dataset for Arabic question answering. https://arxiv.org/abs/2403.17848

[85] Manakul, P., Liusie, A., Mark, J. (2023). SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. arXiv:2303.08896. https://doi.org/10.48550/arxiv.2303.08896

[86] Devlin, J., Chang, M.W., Lee, K., Toutanova, K. (2019). Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 4171-4186. https://doi.org/10.18653/v1/N19-1423

[87] Mubarak, H., Malhas, R., Mansour, W., Mohamed, A., Fawzi, M., Hawasly, M., Elsayed, T., Darwish, K.M., Magdy, W. (2025). IslamicEval 2025: The first shared task of capturing LLMs hallucination in Islamic content. In Proceedings of The Third Arabic Natural Language Processing Conference: Shared Tasks, pp. 480-493. https://doi.org/10.18653/v1/2025.arabicnlp-sharedtasks.67

[88] Tsuruta, H., Sakaguchi, R. (2024). Investigating hallucination tendencies of large language models in Japanese and English. Research Square (Research Square). https://doi.org/10.21203/rs.3.rs-4521710/v1

## A. Translation of the response generated from GPT-4O-MINI

This section provides a translated version of the question and responses using Google Translate as shown in Figure 12.

| | Question | Response Type | Answer |
|---|---|---|---|
| 1 | What are the academic degrees obtained by Muhammad Masoud Al-Awdi? | Main | Mohammed Masoud Al-Awdi is a well-known academic and researcher in his field. He holds several degrees, including: 1. **Bachelor's** in Arabic Language. 2. **Master's** in Arabic Literature. 3. **PhD** in Arabic Literature. If you need more detailed information about his academic career or achievements, feel free to ask more questions! |
| | | Sample 1 | I do not have accurate information about Muhammad Masoud Al-Awdi personally, including details about his academic degrees. You may need to search other sources or official websites to get accurate information about his biography and academics. If you have any other questions or need information about another topic, please feel free to ask! |
| | | Sample 2 | Mohammed Masoud Al-Awdi is a Saudi academic and researcher, who obtained his academic degrees in multiple fields. He obtained a bachelor's degree in chemistry, and a master's degree in chemistry, from King Saud University. He then obtained a doctorate in chemistry from Texas A&M University in the United States |
| | | Sample 3 | Mohammed Masoud Al-Awdi is a prominent academic figure, and has obtained a number of academic degrees. If you are looking for accurate information about his academic degrees, or it is necessary to know his academic affiliations, it is preferable to refer to reliable sources such as official academic websites or his updated CV. Recently, accurate information may not be available to me after October 2021, so you may need to check other sources for the most up-to-date information. |
| 2 | Where was Abdel Wahab Mohamed Abdel Hafeez born? | Main | Abdel Wahab Mohamed Abdel Hafez was born in Dakahlia Governorate, Egypt. If you need more information about him or his achievements, feel free to ask! |
| | | Sample 1 | Abdel Wahab Mohamed Abdel Hafez was born in Minya Governorate, Egypt. If you need more information about him, don't hesitate to ask! |
| | | Sample 2 | Abdel Wahab Mohamed Abdel Hafez is a well-known Egyptian academic and researcher in the fields of humanities and social sciences, but he is not a world-famous figure like artists or athletes. If you have specific context about the person you are asking about, please clarify it so we can better assist you. |
| | | Sample 3 | Abdel Wahab Mohamed Abdel Hafez was born in Kafr El Sheikh, Egypt. If you have any further questions or need additional information, please feel free to ask! |
| 3 | What is the capital of Australia? | Main | The capital of Australia is Canberra. |
| | | Sample 1 | The capital of Australia is Canberra. |
| | | Sample 2 | The capital of Australia is Canberra. |
| | | Sample 3 | The capital of Australia is Canberra. |

**Figure 12.** English translation of example of responses shown in Figure 11.
Note: The Non-factual statements are highlighted in red.