

Evaluation of Quantitative and Qualitative Metrics for Assessing Hallucination Phenomena in Large Language Models

Muhammad Faizan Khan¹

¹The Islamia University of Bahawalpur, Department of Computer Science, Hasilpur, Bahawalpur, Pakistan.

ABSTRACT

Hallucination phenomena in large language models have drawn considerable attention due to their impact on reliability, trustworthiness, and interpretability. Recent advances in transformer-based architectures have demonstrated remarkable capabilities in tasks such as language generation, question answering, and dialogue systems. However, the incidence of fabricated details during inference raises concerns regarding the internal mechanisms that guide these models toward erroneous responses. Uncertainties inherent in model training and data representation create conditions in which hallucinated elements appear, often veiled by fluency and coherence that obscure their inaccuracy. Researchers have proposed numerous strategies to identify and evaluate these outputs, leading to the emergence of a broad array of quantitative and qualitative metrics. Quantitative measurements focus on numerical or probabilistic characterizations, reflecting the extent to which token distributions deviate from reference truths. Qualitative assessments emphasize the interpretive dimension, shedding light on user perceptions, contextual expectations, and semantic coherence. This paper provides a systematic evaluation of these methodologies by examining different metric families, highlighting the conditions under which each approach offers robust insights into generative behavior. The synthesis presented here reveals methodological distinctions, reveals synergies among multiple evaluation frameworks, and showcases promising analytical pathways for furthering the accurate interpretation of model outputs. An integrated perspective on hallucination assessment could guide principled development and deployment of reliable language models.

1 INTRODUCTION

The study of hallucination phenomena in large language models emerged from ongoing efforts to examine discrepancies between generated text and factual content. Researchers working on natural language processing (NLP) have investigated how these models, built upon extensive corpora, can produce fluent outputs that deviate from accurate real-world information. Developer communities observed that, despite improvements in model size and complexity, ungrounded or fabricated details persist as a recurring problem. Scholarly investigations trace the roots of hallucination to model architectures, the distribution of data used in training, and the inherent difficulty of capturing all facets of real-world knowledge in a high-dimensional parameter space.

Transformer-based models, designed to learn vast patterns of language use, illustrate their effectiveness in tasks spanning machine translation, summarization, and conversational agents. Hallucinations, however, raise serious concerns when such models are deployed in contexts involving medical advice, financial guidance, or legal documentation.

The mismatched or invented content may carry real consequences when users rely on system-generated text for decision-making. Researchers have attempted to categorize these hallucinations based on severity, subject domain, and the underlying causes that trigger unfaithful reproduction of facts [1, 2].

Data-driven approaches illuminate the complex interplay between self-attention mechanisms and token distributions that guide generative processes. Large-scale pre-training [3, 4], accompanied by adaptation on specialized datasets [5], has improved the coherence of responses, yet it has not eradicated unintended insertions. Model size increments have been correlated with expansions in representational capacity but have not fully shielded outputs from error. The phenomenon of hallucination appears not only as outright fabrications but also as subtler forms of misalignment where partial truths are blended with invented segments. These observations fuel discussions about the necessity for rigorous evaluation methodologies capable of capturing multi-faceted manifestations of untruthful content.

Category	Description	Example	Impact
Factual Hallucination	Generates incorrect facts	Claiming a wrong historical event date	Misleads users
Inferential Hallucination	Draws incorrect conclusions	Inferring relationships that do not exist	Distorts understanding
Contextual Hallucination	Misinterprets prompt context	Answering an ambiguous query incorrectly	Reduces reliability
Blended Hallucination	Mixes truth with falsehood	Combining real and fake citations	Creates confusion

Table 1. Types of hallucinations in large language models, classified based on characteristics and implications.

Benchmark	Task Type	Hallucination Focus	Limitations
SQuAD	QA	Fact retrieval accuracy	Lacks adversarial cases
TriviaQA	QA	Knowledge consistency	Potential bias in sources
Natural Questions	QA	Grounded response generation	Limited coverage of rare facts
XSum	Summarization	Faithfulness to input	High abstraction leads to errors

Table 2. Common benchmarks used to evaluate hallucinations in NLP models, highlighting their strengths and weaknesses.

Contextual factors contribute to generating fictitious claims that may appear consistent to casual observers. Prompts containing ambiguous or conflicting information can lead the model to fill gaps with invented data. In certain question-answering scenarios, the absence of factual references within the training dataset drives the model to approximate responses, thereby introducing plausible-sounding yet incorrect statements. This challenge becomes more evident in creative tasks, where the boundary between acceptable imaginative expression and factual inaccuracy can blur.

Public repositories of question-answering benchmarks reflect this complexity, offering scenarios where model responses can be partially correct yet garnished with misleading figures or references. Datasets such as SQuAD, TriviaQA, and Natural Questions test a model's retrieval capabilities, but they seldom account for the spontaneous introduction of wholly fabricated items. Newer benchmarks attempt to highlight the phenomenon through curated examples that systematically test the presence of ungrounded details, though these remain evolving areas of investigation. The consistent focus on result accuracy and fidelity to source information underscores the community's need for reliable indicators that distinguish extraneous or incorrect elements in generated text.

Ongoing research shows that hallucinations do not vanish under incremental training or parameter tuning. Instead, they manifest in different ways across tasks that range from narrative generation to summarization. Multi-document summarization efforts provide an illustrative ex-

ample, where references might be conflated or incorrectly merged into a single textual representation, masking the identity of the original sources. The user typically sees a single cohesive summary, unaware that composite hallucinations may have formed at intermediate processing stages. The presence of these distortions has heightened interest in robust scoring methods designed to detect factual discrepancies.

Various approaches have been proposed to measure the degree of fidelity in text generation. Traditional metrics such as BLEU, ROUGE, or METEOR focus on n-gram overlaps with reference texts, yet they fail to capture the logic or truthfulness of the content. Subsequent developments incorporate deeper semantic features, employing knowledge graphs or factual consistency checks to quantify alignment with established facts. These metrics work in tandem with human evaluation protocols that rely on domain experts or crowd-sourced annotators. Even so, the breadth and complexity of language tasks pose significant challenges to ensuring objective and consistent assessment of factual correctness.

Researchers seek to unify these measurement strategies under theoretical frameworks that connect model probability distributions with patterns of factual dependence. Bayesian perspectives, for example, investigate how partial evidence, gleaned from input prompts, influences the posterior probability of each token. Shifts in these probabilities might signal the emergence of fabricated details, offering opportunities for targeted scrutiny. Statistical confidence

Metric	Approach	Strengths	Weaknesses
ROUGE	N-gram Overlap	Measures lexical similarity	Does not check factual accuracy
BLEU	Precision-based Scoring	Useful for structured tasks	Ignores logical correctness
BERTScore	Contextual Similarity	Captures semantic meaning	Sensitive to paraphrasing
FactCC	Fact Consistency Model	Identifies factual deviations	Needs large annotated data

Table 3. Various metrics used to measure hallucinations in text generation, along with their benefits and drawbacks.

intervals, perplexity-based thresholds, and outlier detection algorithms represent only a few of the avenues explored for identifying potential hallucinations. The field lacks consensus on a universal approach, with each methodology offering strengths aligned to specific tasks, data domains, or user requirements.

Technological innovations in model interpretability aim to illuminate how self-attention heads or feed-forward blocks process linguistic cues. Researchers trace the flow of information within the network to identify patterns that correlate with hallucinatory outputs. Visualization techniques attempt to locate segments of text that appear influential in generating spurious claims. These tools facilitate a deeper understanding of the internal dynamics but often remain limited by the scale and complexity of the models. A synergy between interpretability and evaluation might reveal ways to map certain attention patterns onto detectable hallucinations [6, 7].

Academic discourse surrounding hallucination reflects broader questions about accountability in artificial intelligence systems [8]. The assurance that a language model can reliably produce accurate information intersects with ethical discussions about the responsibility of developers and deployers. Interpretive frameworks, shaped by both engineering and philosophical viewpoints, grapple with whether a hallucinating model can ever be deemed trustworthy. This paper adds to these deliberations by systematically presenting a comprehensive landscape of both quantitative and qualitative metrics that aim to capture hallucinatory outputs. Emphasis is placed on methodical rigor rather than on potential remedies or mitigation tactics. Instead, the subsequent sections dissect the core theoretical underpinnings of hallucination, compare current quantitative metrics, explore various qualitative approaches, and consider how an integrated mixed-methods framework might sharpen our ability to detect and interpret such outputs. The ultimate objective is to equip researchers with a robust set of evaluative lenses that accurately reflect the multi-dimensional nature of language model hallucinations [9, 10].

2 FOUNDATIONS OF HALLUCINATION PHENOMENA

Theory-driven perspectives on hallucination phenomena in large language models originate from the tension between probabilistic generative processes and the requirement for factual veracity. Traditional language modeling objectives optimize the likelihood of observed tokens within training corpora. This procedure guides models to generate text that replicates patterns of natural language while implicitly learning contextual associations. However, the learned associations do not uniformly align with real-world truths. The capacity to capture linguistic structure does not guarantee the internalization of consistent factual relationships. Hallucinations arise when the model’s generative mechanism relies on ambiguous or incomplete evidence, resulting in synthetic content that appears coherent but lacks a factual basis.

The architecture of transformers contributes to these phenomena through attention-based weighting of token relationships. Each token interacts with others in a manner that shapes the model’s progressive representation of meaning. The distribution of learned parameters may cluster certain concepts or facts in ways that are not strictly faithful to external reality. Connections formed during the learning process often reflect linguistic co-occurrences more than genuine cause-effect or factual linkages. Such mismatches can lead to overgeneralizations or conflations of separate concepts. In real-world deployments, these conflations translate into text that might seamlessly merge correct data with erroneous additions.

The scaling of parameter counts amplifies the model’s capacity to memorize patterns from vast corpora. Memorization may reduce certain errors by providing broader coverage of factual material. Nevertheless, it also intensifies the risk of generating semantically plausible but factually suspect text. The direct correlation between model size and emergent phenomena, including improved language fluency, underscores why hallucinations remain a persistent concern even in state-of-the-art systems. Researchers suspect that larger models use more complex decision boundaries that can mask mistakes within elaborate contexts.

Probabilistic frameworks emphasize that hallucinations represent an instance of model inference under uncertainty. Bayesian approaches consider the posterior distribution of textual hypotheses given observed input. High-dimensional parameter spaces yield large support for multiple plausible completions. Under constraints such as ambiguous prompts, the model may select from a broad distribution of tokens, increasing the likelihood of speculative or incorrect statements. The interplay between data priors and observation likelihoods influences how these tokens are chosen, resulting in an output that might incorporate partial or incorrect knowledge drawn from the training data.

Language modeling inherently focuses on capturing local consistency rather than global factual coherence. Token-level attention and next-token prediction objectives optimize local transitions, often neglecting the broader factual structure that should govern an entire paragraph or document. This localized optimization can allow for drift from the source truth over the course of multi-sentence generations. Models effectively generate text that matches recognized patterns, yet they do not consistently validate those patterns against real-world facts. This phenomenon illustrates the gap between language modeling objectives and the actual constraints of factual correctness.

Psychological analogies drawn from human cognition describe hallucination as a form of confabulation in which partial memories or associations are combined into new constructs. Such comparisons highlight the model’s predisposition to fill gaps with plausible substitutes, mirroring the human tendency to guess or infer missing details. This viewpoint, though metaphorical, encourages a cross-disciplinary lens on the phenomenon, suggesting that hallucinations may be a natural byproduct of generative processes that rely on probabilistic associations rather than verified knowledge. Researchers harness these analogies to design experiments that test how models respond to uncertain or conflicting cues, revealing points where hallucinations are most likely to emerge.

Empirical investigations into hallucination root causes rely on ablation studies and controlled manipulations of input data. Subsets of training corpora are removed or replaced to observe how the model adapts. Degradation in factual reliability can surface in these scenarios, providing insight into the distribution of parameters responsible for factual recall. Hidden-layer analysis indicates that certain layers or attention heads might be specialized in gleaning factual information, while others focus on linguistic style or other features. The synergy or conflict among these specialized components manifests in varying degrees of hallucination.

Models fine-tuned on specialized data show that alignment with a domain’s factual structure can reduce obvious errors. Such interventions do not eliminate deeper generative tendencies to invent details when faced with knowledge gaps. These findings reinforce the notion that hallucina-

tion cannot be fully attributed to a lack of domain-specific knowledge alone. The generative objective in language modeling prioritizes fluency and local coherence, which can inadvertently overshadow the verification of factual content. The phenomenon thus persists across domains, contexts, and model sizes.

Comparative studies of different architectures, such as recurrent neural networks and convolutional neural networks, indicate that attention-based models often exhibit enhanced representational capacity. That capacity comes with an additional layer of complexity in how information is integrated over long sequences. Hallucinations are still observed in architectures that do not explicitly use the transformer mechanism, implying that the phenomenon is not exclusively tied to one particular design. Instead, the generative approach itself fosters circumstances where plausible but erroneous text is produced, regardless of the specific neural architecture used.

Interpretability research aims to map the internal vectors and attention weights associated with factual recall. Techniques such as layerwise relevance propagation attempt to track how tokens influence the final output. These techniques have uncovered patterns in which words or phrases with minimal initial relevance suddenly influence a segment of text, suggesting that hallucination can sometimes be traced to momentary shifts in attention. Observations of these phenomena underscore the complexity of bridging the gap between a model’s representational depth and transparent human-understandable processes.

The foundations described here set the stage for analyzing quantitative and qualitative metrics designed to evaluate hallucination. A thorough understanding of these phenomena, rooted in theoretical and empirical analyses, aids in interpreting why specific evaluation strategies succeed or fail in detecting fabrications. The subsequent sections present the major families of metrics, illustrating how each addresses various aspects of untruthful model outputs [11, 12].

3 QUANTITATIVE METRICS FOR HALLUCINATION ASSESSMENT

Numerical and algorithmic evaluations of hallucination rely on systematic measurements of content fidelity and consistency. One primary category encompasses reference-based metrics that calculate divergence between generated text and a ground truth. Another category focuses on anomaly detection within a model’s internal probability distributions. The overarching aim is to assign numerical scores that reliably predict the presence or degree of hallucination in a given output [13].

Automatic evaluations originally developed for language tasks, such as BLEU, ROUGE, and METEOR, have been adapted for hallucination detection. BLEU computes n-gram overlaps between candidate and reference sentences, providing a rough gauge of lexical similarity. ROUGE

Metric	Methodology	Strengths	Weaknesses
BLEU	N-gram overlap	Simple and widely used	Ignores factual correctness
ROUGE	Lexical similarity	Effective for summarization	Fails on semantic accuracy
METEOR	Synonym matching	Captures linguistic variations	Limited factual assessment
BERTScore	Embedding similarity	Context-aware comparisons	May overlook factual errors

Table 4. Reference-based metrics commonly used in hallucination assessment, highlighting their methodologies, advantages, and limitations.

compares overlaps of unigrams, bigrams, or longer segments, with particular versions catering to summarization tasks. METEOR incorporates synonym matching and other linguistic features in an effort to capture more nuanced similarity. These traditional metrics do not necessarily capture factual correctness, since linguistic similarity does not equate to factual validity. Researchers have attempted to refine these scores by including synonyms relevant to domain knowledge or weighting certain terms more heavily if they refer to verifiable entities, but the fundamental shortcoming remains that an output can share lexical patterns with a reference while introducing subtle factual errors.

Embedding-based metrics such as BERTScore address some limitations by comparing sentence embeddings derived from large pretrained language models. Textual similarity is measured at a contextualized token level, capturing semantic equivalences beyond mere string overlap. BERTScore aligns tokens in the candidate and reference based on similarity in the embedding space. This provides a more flexible method for gauging semantic closeness. However, it still does not fully capture erroneous facts, since two passages might share a high-level semantic topic but differ in factual detail. Large models used for embedding generation could themselves harbor incomplete or inaccurate knowledge, which complicates the reliability of embedding-based assessments.

Factual consistency metrics attempt a targeted approach by checking the presence or correctness of specific facts. FactCC, for instance, applies a binary classification model that takes pairs of text segments—one from the source document and one from the summary—to predict if the summary is consistent with the source. This method uses data automatically generated by introducing known factual distortions in text to train consistency classifiers. Another approach leverages knowledge graphs, extracting relational triples from generated text and comparing them against established databases or structured knowledge. Discrepancies in subject-predicate-object relationships can indicate hallucination. These techniques reduce reliance on broad lexical similarity and instead focus on explicit fact verification.

Perplexity-based methods analyze how probable a given

output is under a secondary or reference model. Low perplexity suggests that an output aligns well with learned distributions of text, yet it does not necessarily correspond to factual veracity. To bridge this gap, some work examines perplexity shifts when critical factual tokens appear. Sudden increases or decreases in perplexity around named entities, dates, or other verifiable details may signal attempts by the model to produce material outside of its core knowledge. Such signals, when aggregated, can help detect areas where hallucination is more likely. The success of perplexity-based detection hinges on choosing an appropriate reference model or distribution, which remains a nontrivial design choice.

Outlier detection frameworks formulate hallucination identification as an anomaly detection task. Generated tokens that deviate significantly from expected embeddings or distributional properties are flagged for scrutiny. Statistical modeling of token embeddings can highlight segments of text with unusual patterns, suggesting that the model ventured outside the typical scope of its learned parameters. Clustering methods can isolate groups of tokens or phrases that consistently correlate with erroneous statements. This approach reduces the reliance on external references, instead probing the self-consistency of the model’s internal representation.

Scoring systems that quantify source-reference alignment have grown in complexity to accommodate multi-document inputs or partial references. Summaries derived from multiple sources are compared against each source to check for contradictions or distortions. Weighted scoring aggregates these comparisons into a final measure of how faithfully the generated text matches the ensemble of original information. The weighting function may account for the reliability of each source. Instances in which the model incorrectly merges or synthesizes incompatible details are penalized, highlighting occurrences of hallucination. Multi-document scenarios underscore the need for robust alignment metrics that navigate conflicting or overlapping facts.

Many quantitative metrics are evaluated through correlation studies with human annotations [14, 15]. Researchers

Method	Technique	Hallucination Detection	Challenges
FactCC	Binary classification	Identifies factual inconsistencies	Requires annotated data
Knowledge Graphs	Relation extraction	Matches facts to databases	Limited by knowledge coverage
Perplexity Shifts	Probability distribution	Detects uncertain token generations	Choice of reference model matters
Entailment Models	Logical inference	Evaluates consistency with source text	Sensitive to paraphrasing

Table 5. Key factual consistency evaluation techniques, their detection capabilities, and associated challenges.

Metric	Combination Approach	Advantages	Limitations
Hybrid Scoring	Embeddings + factual checks	Balances lexical and factual accuracy	Requires complex integration
Multi-Document Alignment	Cross-referencing multiple sources	Detects synthesis hallucinations	Hard to handle conflicting facts
Outlier Detection	Embedding anomaly detection	Identifies deviations from norm	Needs fine-tuned thresholds
Weighted Fact Matching	Entity verification with weighting	Penalizes severe factual errors	Domain-specific tuning required

Table 6. Hybrid and advanced hallucination detection metrics, illustrating their methodologies, advantages, and constraints.

assemble labeled datasets where domain experts or crowd workers identify hallucinated segments [16]. The correlation between metric outputs and these human judgments serves as a key indicator of effectiveness. High correlation suggests the metric can replicate human detection of invented content. Cases with low correlation often reveal that the metric either overemphasizes superficial language features or fails to capture subtle factual discrepancies. Studies using correlation typically note the importance of broad coverage across multiple domains and tasks to ensure generalizability.

Hybrid metrics combine the best features of the approaches described above. A system might first apply an embedding-based comparison to identify semantically aligned passages, then run a factual consistency check on critical entities. Weighting schemes account for the severity of discrepancies: minor factual slips might be scored less harshly than gross inventions. Some research incorporates textual entailment models that assess if the generated text is logically entailed by source documents. Entailment-based scores attempt to unify lexical, semantic, and factual checks into a cohesive measure. The success of these integrated methods depends on carefully curated training data and the design of robust classification or regression models.

While many quantitative metrics show promise in flagging certain classes of hallucination, they are not universally reliable. Strategic analysis of failure cases shows that metrics can be fooled by text that preserves key entities but alters crucial relationships. Overly rigid reference-based mea-

sures may penalize legitimate paraphrases or expansions that remain factually correct. Overly flexible embedding-based measures may miss subtle factual inversions if the semantic context remains largely unchanged. The interplay between these factors illustrates the balancing act required to build a single numeric measure that captures the essence of hallucination.

Scientific communities continue to refine evaluation datasets dedicated to hallucination detection. Efforts to construct synthetic corpora that systematically introduce distortions provide valuable training and benchmarking tools for new metrics. Real-world corpora annotated at a granular level enable comparative analyses of how each approach handles domain-specific nuances. These initiatives add depth to the field, revealing new techniques for better modeling the intricacies of generated errors.

4 QUALITATIVE EVALUATIONS OF HALUCINATION

Subjective interpretations and user-centric approaches bring essential nuance to the assessment of hallucination in large language models. Qualitative evaluations aim to capture dimensions that purely numerical metrics may overlook. User judgment, contextual relevance, and interaction patterns all shape how a fabricated detail is perceived. These methods involve systematic protocols that parse output text for logic, coherence, and consistency with domain knowledge, but rely on direct human engagement to interpret the severity

and nature of errors.

Expert reviews represent a cornerstone of qualitative evaluation, deploying domain specialists to critique model outputs. Medical or legal experts, for instance, can pinpoint discrepancies that a layperson would not recognize. Judgments from such assessments yield insights into how a model's invention of details correlates with domain-specific standards of accuracy. The specialized knowledge of these experts enables them to identify subtle or deeply embedded inconsistencies that are not easily flagged by simple text matching. Qualitative feedback from expert reviews often informs guidelines for refining or recalibrating models, although that facet of engineering is outside the scope of this discussion.

Crowdsourced annotations form a second major technique, capitalizing on the diverse perspectives of lay annotators. Platforms that host large pools of workers facilitate rapid evaluation at scale, often with multiple annotators assessing each output. These workers may receive task-specific instructions emphasizing factual correctness, consistency, and clarity. They label portions of text deemed suspect, unclear, or overtly incorrect. Majority voting and inter-annotator agreement scores gauge the reliability of these human judgments. Although crowdsourced evaluations lack the specialized depth of expert assessments, the broad demographic representation can highlight how typical end-users might perceive or react to hallucinated content in more general contexts.

Holistic rating scales, known from fields like qualitative research in psychology and user experience studies, have been adapted for language model assessment. Annotators read entire passages or dialogues and rate them on scales that encompass fluency, semantic coherence, and factual correctness. Low ratings on factual correctness often coincide with identified hallucinations. This approach attempts to capture the integrated experience of encountering model-generated text, reflecting not just snippet-level issues but also broader contexts and narrative flow.

Protocol analysis techniques involve asking annotators to articulate their reasoning when identifying hallucinated material. This is accomplished through think-aloud methods or structured interviews in which participants explain their thought process while reading the output. Researchers transcribe these discussions, subsequently coding them for recurring themes such as confusion points, triggers for suspicion, or verification processes. In domain-specific tasks, participants may reference external resources or domain expertise as they verify the text. This granular viewpoint helps clarify which aspects of an output appear the most deceptive or plausible, thereby informing future improvements in detection strategies.

Discourse analysis provides another qualitative lens, examining how the model constructs narrative or argumentation. Linguistic features such as discourse markers, coherence relations, and thematic progression indicate whether

the text remains internally consistent. Sudden shifts in topic or contradictory statements can signal potential hallucination. Text that reads as logically incoherent might not always be flagged by purely quantitative methods. Researchers trained in discourse analysis investigate how the structure of the text might cue or obscure recognition of inaccuracies.

Contextualized scenario testing offers a methodology in which users or annotators interact with the model in real time. The model's hallucinations are identified based on interactive questioning, clarifications, or contradictory follow-up prompts. This approach is relevant for applications like customer service chatbots or educational tutoring systems, where the user continuously evaluates the system's consistency. Hallucinations can emerge through dynamic interactions, and the user's immediate response might trigger corrections, leading to iterative cycles of revelation and resolution. Although the correction strategies are not under discussion here, the recognition of the phenomenon itself is central to understanding the model's quality.

Some qualitative evaluations employ rhetorical analysis, focusing on persuasive or emotive tactics embedded in the text. In certain generative tasks, the model may produce content that is factually incorrect but employs rhetorical devices to appear convincing. Annotators trained in argumentation analysis document how these rhetorical strategies may cloak factual errors. These strategies might be especially pronounced when the model attempts to fill knowledge gaps through authoritative-sounding language, referencing non-existent sources or misquoting real sources. Qualitative inspections help illustrate how textual style and structure can serve as vectors for deceptive impressions.

Cultural and societal dimensions often surface in qualitative evaluations, recognizing that what constitutes a hallucination might vary across different communities or languages. Misrepresenting historical events or cultural practices can be considered severely incorrect in one context but might be less scrutinized in another. Cross-cultural annotation efforts involve multi-lingual or multi-cultural annotators verifying the text against diverse sets of knowledge. Researchers observe that certain language models, trained predominantly on data from specific regions, might produce hallucinatory statements when addressing less familiar cultures or languages. Qualitative annotations thus highlight the importance of context when diagnosing generative inaccuracies.

Synthesis of qualitative findings often leads to thematic coding that groups hallucination incidents by type. Overlapping categories may emerge, such as conflation of separate entities, fabrication of references, or misalignment of numerical data. These coded patterns then inform the design of specialized detection protocols or the adaptation of domain-relevant checklists. Although these designs pertain to improvements in detection, the raw data and annotated feedback themselves provide valuable empirical founda-

tions for understanding how hallucinations materialize in real usage scenarios.

Human judgment inevitably carries variability and subjective bias. Inter-annotator agreement studies are used to monitor this variability, employing statistical measures like Cohen's kappa or Krippendorff's alpha to quantify consistency among raters. Results that show high agreement suggest the presence of hallucinated material that is recognizable without extensive dispute. Low agreement might indicate borderline cases or subtle forms of misinformation that only some annotators detect. Such findings guide further refinement of annotation guidelines, ensuring a stable qualitative framework.

Comprehensive comparisons of qualitative evaluations with quantitative metrics reveal that each approach captures distinct aspects of hallucination. Numerical scores excel at scale and objectivity, while human-based judgments interpret complex or contextualized inaccuracies. When used in conjunction, these perspectives can provide a multi-dimensional portrait of how and why hallucinations occur, as well as the extent to which they might influence user trust or comprehension. The next section discusses the integrated application of both quantitative and qualitative approaches, elucidating strategies that leverage mixed methods to offer a fuller evaluation landscape for hallucination phenomena.

5 MIXED-METHODS INTEGRATION IN HALUCINATION ANALYSIS

Evaluation frameworks that incorporate both quantitative and qualitative metrics offer a more complete picture of hallucinatory tendencies in large language models. Hybrid methodologies leverage the strengths of numeric objectivity while retaining the contextual nuance that human judgments provide. Researchers aim to integrate these two assessment paradigms through coordinated data collection, result comparison, and iterative refinement, resulting in robust insights into the behavior of generative models.

Multi-phase evaluation pipelines often begin with a broad quantitative screening of outputs to identify potentially hallucinated examples. Automated systems score each example according to one or more numeric metrics, ranking them by likelihood of containing errors. Subsequent steps enlist human annotators to thoroughly review a selected subset of examples, focusing on those deemed most suspicious by the numeric scores, along with a smaller sample randomly chosen from the entire distribution. This design helps balance efficiency with coverage, ensuring that systematic biases in the automatic metrics do not overshadow the qualitative appraisal.

Interactive tools facilitate real-time triangulation between human opinions and metric outputs. Researchers might deploy graphical interfaces where annotators can see numeric scores for each sentence or phrase within a generated text. These scores often incorporate perplexity shifts, factual consistency checks, or anomaly detection

flags. Users can then confirm, refine, or contradict the automatic assessments, providing detailed remarks on where and why the metric may have misjudged a passage. This process harnesses the machine's speed in sifting through large corpora and the human's capacity for contextual discernment.

Iterative improvement cycles form a core component of mixed-methods strategies. Investigators compile instances where numeric metrics and human judgments diverge, analyzing them to identify possible patterns or root causes. In some cases, the metric might penalize unusual but factually correct references due to incomplete training or flawed assumption sets. In others, the metric might overlook embedded misinformation that humans readily identify. Aggregated findings guide modifications to the metric's weighting or the algorithm's classification thresholds. The updated system undergoes further evaluation in a continuous feedback loop with human reviewers, driving incremental enhancement of metric reliability.

Semantic mapping approaches seek to formalize the relationship between numeric scores and conceptual categories identified by qualitative annotations. A set of hallucinatory categories, for instance, might be defined to capture repeated error types, such as reference fabrication, conflation of related entities, or duplication of partial truths. Statistical modeling ties these categories to distributions of metric values. Researchers can then forecast the likelihood that a certain range of metric scores corresponds to a specific category of hallucination. Human-labeled data is essential for this modeling, as it trains the system to recognize which numeric patterns correlate with domain-relevant error types.

Quantitative metrics can be used as a pre-screening step for deeper, domain-focused analyses. In complex fields like biomedical text generation, an automated tool might flag outputs with potential factual inconsistencies, while specialized annotators perform a deeper inspection of the flagged sections. Domain experts evaluate each flagged segment against scientific literature or medical guidelines, determining whether it constitutes a hallucination in context. The synergy between broad computational scanning and targeted expert appraisal maximizes efficiency when addressing large-scale corpora.

Correlation coefficients between various metrics and human judgments provide an overarching measure of how well different quantitative approaches align with subjective perceptions of accuracy. Mixed-methods research often aims to exceed moderate correlation thresholds, such as Spearman's or Pearson's values in the range of 0.70 or above. Achieving such figures consistently across multiple tasks and domains offers evidence that the proposed metric resonates with human evaluations of factual correctness. If correlations falter in certain domains, refined domain-specific metrics may be introduced or domain experts may be included in the review process to increase granularity.

Advanced annotation platforms sometimes incorporate model feedback loops, where a model trained to detect hallucinations is updated based on both numeric and qualitative indicators. The model iteratively adjusts its parameters to match the consensus derived from these combined signals. Over time, the hallucination-detection system itself gains the capacity to approximate human judgments more closely, reaching higher accuracy levels than standalone numeric or purely human-based methods. Such systems must continually incorporate new data and new types of error to maintain relevance in dynamic deployment environments.

Meta-analyses of published findings serve a higher-level function, synthesizing results from multiple studies that employ different combinations of metrics. By pooling data, researchers can ascertain which metrics maintain robust performance across tasks, languages, or content domains. The meta-analytic perspective clarifies whether certain approaches generalize widely or if they are limited to specialized scenarios. It also identifies missing links in the field's collective knowledge base, indicating areas for further methodological innovation.

Case-based synergy between quantitative outputs and qualitative feedback informs how end-users experience hallucinatory text. Although not delving into case studies as a formal section, contextual examples help illustrate how numeric scores might flag a summary as suspicious due to unusually low alignment with the source, and how a human review might confirm that the text indeed introduced non-existent results or references. Detailed dialogue between these two modalities fosters confidence in the final determination.

Surveys of actual user populations, alongside demonstration sessions, can reveal points of alignment or mismatch between user impressions and the integrated metrics. Some users might prioritize coherence over factual precision, viewing minor inaccuracies as acceptable in creative tasks. Others might demand strict factual adherence in critical applications. The integrated framework adapts to these expectations by calibrating metrics and annotation schemes that reflect real-world usage priorities. This calibration can lead to refined weighting schemes that highlight factual correctness more strongly in crucial domains like finance or healthcare, while granting more creative leeway in entertainment applications.

Hybrid evaluation results are typically reported in a multi-layered format, providing numeric scores for each metric, integrated scores that combine them, and descriptive commentary or coded annotations from human reviewers. The layered approach offers a fuller account of where the model succeeds or fails, granting researchers and stakeholders a transparent view of the system's performance. Quantitative ranks and thresholds provide swift comparability between different versions or models, and the qualitative narrative unpacks subtle errors that might not be captured in a single number.

Reliance on mixed methods addresses the intricacies of hallucination by acknowledging that a single vantage point cannot fully characterize the phenomenon. Automated scoring remains valuable for scale and speed, while human assessment brings the depth of contextual knowledge and interpretation. Their combined application generates a versatile toolkit for robust evaluation, forming the basis for future refinements in the objective measurement and nuanced understanding of model-generated inaccuracies [17, 18].

6 CONCLUSION

Evaluations of hallucination phenomena in large language models benefit from a multifaceted approach that incorporates both quantitative and qualitative metrics. The quantitative lens, exemplified by reference-based overlap measures, factual consistency checks, perplexity analyses, and anomaly detection, provides scalable and objective indicators of unfaithful output. These numeric systems enable rapid screening across large datasets [19], but they often fail to capture subtle or context-specific errors that might appear plausible to automated scoring methods. Qualitative frameworks, supported by expert reviews, crowdsourced annotations, protocol analyses, and discourse studies, address these gaps by uncovering more nuanced manifestations of inaccuracy. Their reliance on human insight yields deeper contextual knowledge and illuminates how errors affect user perceptions of trust and understanding [20, 21].

Synthesis of these methodologies underscores the value of mixed-methods integration. Iterative processes that merge automated scoring with targeted human examination allow for refined metric calibrations and consistent improvements in detection sensitivity. Consensus-building among various metric outputs, complemented by expert commentary, enriches the reliability of final evaluations. This integrated perspective fosters a detailed understanding of the circumstances under which hallucinations proliferate, advancing the design of more precise evaluative instruments. Future developments in assessment are likely to explore more granular correlations between numeric scores and thematically coded qualitative findings, broadening the capacity for accurate measurement across diverse tasks and domains.

Comprehensive methodology consolidations highlight the complexity of generative text alignment with real-world truths. The phenomenon of hallucination persists despite larger model sizes and seemingly more sophisticated training regimens. Ongoing studies indicate that attention-based architectures, probability distributions, and domain-specific factors interact in ways that permit factual inventions to surface. The absence of universal agreement on a single metric or approach further confirms that a diverse suite of evaluations is essential to capture the full spectrum of hallucinatory behavior. Continuous investment in this research area is expected to yield increasingly nuanced tools that enhance our understanding and detection of these phenomena in practical implementations.

REFERENCES

- [1] Vu, T. *et al.* Freshllms: Refreshing large language models with search engine augmentation. *arXiv preprint arXiv:2310.03214* (2023).
- [2] Rawte, V., Sheth, A. & Das, A. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922* (2023).
- [3] Azamfirei, R., Kudchadkar, S. R. & Fackler, J. Large language models and the perils of their hallucinations. *Critical Care* **27**, 120 (2023).
- [4] Zhang, Y. *et al.* Siren's song in the ai ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219* (2023).
- [5] Bhaskaran, S. V. Enterprise data architectures into a unified and secure platform: Strategies for redundancy mitigation and optimized access governance. *Int. J. Adv. Cybersecurity Syst. Technol. Appl.* **3**, 1–15 (2019).
- [6] Ouyang, Q., Wang, S. & Wang, B. Enhancing accuracy in large language models through dynamic real-time information injection. (2023).
- [7] Chen, B., Zhang, Z., Langrené, N. & Zhu, S. Unleashing the potential of prompt engineering in large language models: a comprehensive review. *arXiv preprint arXiv:2310.14735* (2023).
- [8] Mehta, R., Hoblitzell, A., O'keefe, J., Jang, H. & Varma, V. Halu-nlp at semeval-2024 task 6: Metacheckgpt-a multi-task hallucination detection using llm uncertainty and meta-models. In *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, 342–348 (2024).
- [9] Elaraby, M. *et al.* Halo: Estimation and reduction of hallucinations in open-source weak large language models. *arXiv preprint arXiv:2308.11764* (2023).
- [10] Dziri, N., Milton, S., Yu, M., Zaiane, O. & Reddy, S. On the origin of hallucinations in conversational models: Is it the datasets or the models? *arXiv preprint arXiv:2204.07931* (2022).
- [11] Niu, C. *et al.* Ragtruth: A hallucination corpus for developing trustworthy retrieval-augmented language models. *arXiv preprint arXiv:2401.00396* (2023).
- [12] Dhuliawala, S. *et al.* Chain-of-verification reduces hallucination in large language models. *arXiv preprint arXiv:2309.11495* (2023).
- [13] Mehta, R., Hoblitzell, A., O'Keefe, J., Jang, H. & Varma, V. Metacheckgpt—a multi-task hallucination detection using llm uncertainty and meta-models. *arXiv preprint arXiv:2404.06948* (2024).
- [14] Mündler, N., He, J., Jenko, S. & Vechev, M. Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation. *arXiv preprint arXiv:2305.15852* (2023).
- [15] Liang, X. *et al.* Uhgeval: Benchmarking the hallucination of chinese large language models via unconstrained generation. *arXiv preprint arXiv:2311.15296* (2023).
- [16] Bhaskaran, S. V. A comparative analysis of batch, real-time, stream processing, and lambda architecture for modern analytics workloads. *Appl. Res. Artif. Intell. Cloud Comput.* **2**, 57–70 (2019).
- [17] Li, J., Cheng, X., Zhao, W. X., Nie, J.-Y. & Wen, J.-R. Halueval: A large-scale hallucination evaluation benchmark for large language models. *arXiv preprint arXiv:2305.11747* (2023).
- [18] Ji, Z. *et al.* Survey of hallucination in natural language generation. *ACM Comput. Surv.* **55**, 1–38 (2023).
- [19] Bhaskaran, S. V. Tracing coarse-grained and fine-grained data lineage in data lakes: Automated capture, modeling, storage, and visualization. *Int. J. Appl. Mach. Learn. Comput. Intell.* **11**, 56–77 (2021).
- [20] Huang, L. *et al.* A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232* (2023).
- [21] Guan, T. *et al.* Hallusionbench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. *arXiv preprint arXiv:2310.14566* (2023).