# Improving Safety and Reliability in a Conversational AI Assistant

## 1. Introduction

Conversational AI systems are increasingly deployed in real-world settings where users rely on them for information, guidance, and decision support. As these systems scale, issues related to reliability and safety become more pronounced. In this document, I analyze four observed problems in a conversational AI assistant—inconsistency across turns, hallucination, demographic bias, and prompt sensitivity—and propose technically grounded solutions. I focus especially on hallucination and inconsistency, as these most directly undermine user trust and safety.

## 2.1 Problem Analysis

### Inconsistent Responses Across Turns

Inconsistency occurs when the model contradicts its own previous statements within the same conversation. The most likely underlying cause is that the model does not maintain a structured or grounded representation of conversational state. While transformer-based models condition on recent context tokens, they do not explicitly track commitments or factual assertions made earlier. In addition, decoding stochasticity and temperature-based sampling can introduce variability across turns, even when the semantic intent remains unchanged.

To measure inconsistency, I would use a combination of automated and human evaluation. One automated approach is to treat earlier model statements as premises and later responses as hypotheses, then apply a natural language inference (NLI) model to detect contradictions. Metrics such as "contradiction rate per conversation" or "self-contradiction frequency over N turns" would provide quantitative signals. Human evaluation is also important, particularly for nuanced or implicit contradictions.

### Hallucination

Hallucination refers to the model generating factually incorrect or fabricated information with high confidence. This is largely a consequence of the next-token prediction objective used during training, which rewards fluency and plausibility rather than factual correctness. When the model lacks sufficient knowledge or evidence, it often fills gaps by generating statistically likely continuations rather than expressing uncertainty. The absence of grounding mechanisms, such as retrieval or verification, further exacerbates this issue.

Hallucination can be measured using factuality benchmarks such as TruthfulQA or domain-specific QA datasets where ground truth is available. Another useful metric is the

"supportedness rate," which measures whether generated factual claims can be verified against trusted sources. Calibration metrics, such as expected calibration error (ECE), are also useful to identify cases where the model is overconfident but incorrect.

## Bias

Bias arises when the model produces systematically different or harmful outputs based on demographic attributes such as gender, race, or religion. These behaviors typically originate from imbalances and stereotypes present in the training data, which the model absorbs and amplifies. Bias can also be reinforced during fine-tuning if feedback signals are unevenly distributed across demographic contexts.

Bias can be quantified using controlled prompt templates where only demographic attributes are varied, as well as benchmark datasets such as StereoSet or CrowS-Pairs. Metrics might include differential toxicity scores, sentiment shifts, or preference asymmetries across demographic groups.

## Prompt Sensitivity

Prompt sensitivity occurs when small changes in phrasing lead to disproportionately different outputs. This often reflects shallow pattern matching, underspecified objectives, and a lack of robustness to paraphrase. It may also be amplified by instruction-tuning procedures that overfit to specific prompt formats.

Prompt sensitivity can be measured by constructing paraphrase sets for the same intent and evaluating output variance. Metrics such as answer consistency scores or semantic similarity across outputs provide useful signals.

## Prioritization

While all four issues are important, I would prioritize **hallucination** and **inconsistency** first. These issues most directly undermine user trust and can lead to real-world harm, especially in informational or advisory settings. Addressing hallucination and inconsistency also tends to reduce prompt sensitivity as a secondary effect, since grounded and state-aware systems are more stable overall.

# 2.2 Proposed Solutions

## Solution 1: Reducing Hallucination Through Grounding and Verification

To reduce hallucination, I would introduce a grounding-based approach that combines retrieval augmentation with post-generation verification. For queries involving factual information, the model would retrieve relevant documents from a trusted corpus and condition its response on that evidence. Generated answers would include citations or references to retrieved sources. Additionally, a lightweight verification step could extract factual claims from the response and check them against retrieved documents using an entailment model.

This approach would require access to a curated knowledge corpus, a retrieval system (e.g., dense embeddings with approximate nearest neighbor search), and additional compute for retrieval and verification. A realistic timeline for a prototype system would be 6–8 weeks, including data preparation, integration, and evaluation.

Success would be evaluated using factuality benchmarks, supportedness rates, and reductions in high-confidence incorrect answers. Calibration improvements would also be a key metric. Potential risks include increased latency, failure cases when retrieval returns poor results, and over-refusal behavior if the system becomes overly conservative.

**Solution 2: Improving Consistency with Explicit Conversational State**

To address inconsistency, I would introduce an explicit representation of conversational state that tracks factual commitments made by the model. After each turn, key assertions could be extracted and stored in a structured memory. During subsequent turns, this memory would be retrieved and incorporated into the model's context, either as additional conditioning text or through a separate memory-attention mechanism.

Training could be augmented with consistency-aware objectives, such as penalizing contradictions detected by an NLI model. Preference-based fine-tuning could also be used, where consistent responses are ranked higher than contradictory ones.

This solution would require annotated conversational data or synthetic contradiction pairs, moderate additional compute for training, and careful design to avoid reinforcing incorrect earlier statements. Metrics would include contradiction rates, entity consistency, and human judgments of coherence. Risks include error propagation if incorrect facts are stored in memory and reduced flexibility when legitimate corrections are needed.

# 2.3 Experimental Design

To evaluate the hallucination-reduction approach, I would design an experiment comparing a baseline conversational assistant against a retrieval-augmented version with verification.

The primary hypothesis is that grounding and verification will significantly reduce unsupported factual claims without substantially degrading helpfulness. The control condition would be the existing model, while the treatment condition would include retrieval and citation mechanisms.

The dataset would consist of a diverse set of factual user queries, including both answerable and ambiguous questions. A sample size of several hundred queries would allow for statistically meaningful comparisons. Each response would be evaluated for factual correctness, supportedness, confidence, and helpfulness, using a mix of automated checks and human ratings.

Statistical analysis would involve paired comparisons between control and treatment responses, with confidence intervals and significance testing on supportedness and correctness metrics. If supportedness improves while helpfulness remains stable or only

slightly decreases, the intervention would be considered successful. A large drop in helpfulness would indicate the need to adjust refusal thresholds or retrieval quality.

## 2.4 Broader Implications

Introducing grounding and consistency mechanisms will inevitably affect model capabilities. While factual reliability improves, response latency may increase, and the model may appear less fluent or creative in some contexts. There is an inherent trade-off between safety and unconstrained generative performance.

Communicating these changes to users is important. Users should understand when responses are grounded in sources, when uncertainty is present, and why the model may refuse to answer certain questions. Clear explanations and transparent design can help maintain trust even when the model is more cautious.

Ultimately, prioritizing safety and reliability leads to systems that users can depend on over the long term, even if it requires accepting some constraints on raw generative flexibility.