# Improving Safety and Reliability of Conversational AI Systems

## 2.1 Problem Analysis

### 2.1.1 Inconsistent Responses Across Turns

Likely Causes

1) Limited long-term conversational state: Transformer-based models rely on finite context windows and implicit state tracking, leading to drift in beliefs over turns.
2) Probabilistic decoding: Stochastic sampling (e.g., temperature, top-p) introduces variability that can surface contradictions.
3) Lack of explicit world-model constraints: The model does not enforce global consistency across generated statements.

Measurement & Quantification

1) Contradiction Rate: Use Natural Language Inference (NLI) models to detect contradictions between earlier and later turns.
2) Self-Consistency Score: Ask the model the same factual questions at different turns and measure answer agreement.
3) Conversation-Level Consistency Benchmarks (e.g., TruthfulQA-style multi-turn extensions).

### 2.1.2 Hallucination (Fabricated Facts)

Likely Causes

1)  Next-token prediction objective: The model optimizes fluency over factual grounding.
2) Sparse or outdated training data: Missing knowledge encourages plausible-sounding fabrication
3) Overconfidence from RLHF: Human preference training may reward confident-sounding answers.

Measurement & Quantification

1) Factual Error Rate (FER): Percentage of generated factual claims that are incorrect, evaluated against trusted knowledge bases.
2) Attribution Accuracy: Ability to correctly cite sources when prompted.
3) Truthfulness Benchmarks: TruthfulQA, FEVER-style claim verification.

### 2.1.3 Demographic Bias

Likely Causes

1) Bias in pretraining data: Internet-scale corpora reflect societal biases.

2) Spurious correlations: The model learns shortcuts associating demographics with outcomes.
3) Insufficient counterfactual data: Lack of balanced examples during training.

Measurement & Quantification

1) Group Fairness Metrics: Differences in sentiment, toxicity, or refusal rates across demographic groups.
2) Bias Benchmarks: StereoSet, CrowS-Pairs.
3) Counterfactual Evaluation: Swap demographic attributes and measure output changes.

## 2.1.4 Prompt Sensitivity

Likely Causes

1) Highly non-linear decision boundaries in embedding space.
2) Instruction-following overfitting: Sensitivity amplified by RLHF.
3) Lack of robustness objectives during training.

Measurement & Quantification

1) Output Variance under Paraphrasing: Measure semantic divergence for paraphrased prompts.
2) Robustness Curves: Performance vs. degree of prompt perturbation.

## 2.1.5 Prioritization

Top Priorities:

1) Hallucination – Directly impacts trust and safety, especially in medical, legal, and educational contexts.
2) Inconsistent Responses – Undermines reliability in multi-turn interactions and agentic workflows.

Bias and prompt sensitivity are critical but can be addressed iteratively alongside these core reliability issues.

# 2.2 Proposed Solutions

## Priority 1: Hallucination Reduction via Retrieval-Augmented Generation (RAG) + Uncertainty Modeling

Technical Approach

1) Integrate a retrieval module (e.g., dense vector search) to fetch relevant documents.
2) Condition generation on the retrieved evidence.
3) Add selective generation with abstention: when retrieval confidence is low, the model responds with uncertainty.
4) Fine-tune with factuality-aware loss, penalizing unsupported claims.

Required Resources

1) Data: Curated knowledge base, fact-checked QA pairs.
2) Compute: Moderate GPU resources for fine-tuning; CPU-heavy retrieval infra.
3) Timeline: 6–8 weeks (infra + fine-tuning + evaluation).

Evaluation Metrics

1) Factual Error Rate
2) Answerable vs. Abstained Accuracy
3) User Trust Scores (human eval)

Risks & Limitations

1) Retrieval latency
2) Knowledge base coverage gaps
3) Over-abstention reducing usefulness


## Priority 2: Consistency via Memory-Augmented and Self-Verification Techniques

Technical Approach

1)  Introduce an explicit conversation memory storing key facts and commitments.
2)  Use self-verification loops: generate → critique → revise.
3)  Apply consistency regularization during fine-tuning using synthetic contradiction data.

Required Resources

1) Data: Multi-turn dialogues with annotated contradictions.
2) Compute: Additional inference-time cost for verification passes.
3) Timeline: 4-6 weeks

Evaluation Metrics

1) Contradiction Rate
2) Multi-turn QA Accuracy
3) Human consistently ratings

Risks & Limitations

1) Increased inference cost
2) Possible over-constraining of creative responses

# 2.3 Experimental Design

## Experiment: Evaluating RAG for Hallucination Reduction

Hypothesis Retrieval-augmented generation significantly reduces hallucination without degrading answer usefulness.

Experimental Setup

1) Control: Base LLM (no retrieval).
2) Treatment: LLM + RAG + abstention mechanism.

Data & Sample Size

1) 1,000 factual QA prompts across domains (science, history, medicine).
2) Power analysis targeting detection of ≥5% FER reduction.

Statistical Analysis

1) Paired t-test or bootstrap CI on FER differences.
2) Secondary analysis on abstention rates and human usefulness scores.

Expected Outcomes

1) FER significantly lower in treatment.
2) Slight increase in abstentions, acceptable if usefulness remains stable.

Interpretation

1) If FER ↓ and usefulness ↔: successful.
2) If FER ↓ but usefulness ↓: tune abstention threshold.
3) If no FER change: investigate retrieval quality.

# 2.4 Broader Implications

Impact on Model Capabilities

1) Improved trustworthiness and deployment readiness.
2) Slightly reduced fluency or creativity in edge cases.

Safety vs. Performance Trade-offs

1) Higher latency and compute cost.
2) Conservative responses may frustrate some users.

User Communication

1) Transparently communicate uncertainty handling (e.g., "I may be mistaken").
2) Provide citations and explain when the model abstains.

# References

Lewis et al., *Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks*

Lin et al., *TruthfulQA*

Manakul et al., *Self-Verification Reduces Hallucination in LLMs*

Zhao et al., *Calibrating Language Models*