

# Ambiguous by Design: Testing LLMs on Polysemy, Homonymy, and Underspecification

## 1. Abstract

Large Language Models (LLMs) continue to show remarkable fluency & versatility across a range of tasks. However, their performance is heavily influenced by the structure and clarity of the prompts they receive. This study investigates a critical but underexplored axis of prompt design: lexical ambiguity. We construct a controlled dataset of prompt pairs, each containing an ambiguous query and its corresponding disambiguated version, and evaluate outputs from multiple state-of-the-art LLMs (GPT-4, Claude 4 Sonnet, and Gemini 2.5 Pro). Using a combination of semantic similarity (BERTScore), hallucination detection (NER-based fact-checking), and confidence calibration metrics (hedging frequency, assertiveness index), we quantify the effects of ambiguity on output quality. Our results show that ambiguous prompts significantly increase the likelihood of misinterpretation and factual hallucination, and often reduce the model's confidence calibration. Moreover, we observe model-specific variance in ambiguity robustness, suggesting prompt sensitivity is architecture-dependent. These findings illuminate a critical failure point in prompt-response alignment and offer empirical insights for improving prompt engineering practices. This work represents one of the first empirical investigations of lexical ambiguity in LLM interaction, with implications for prompt safety, automated reasoning, and model interpretability.

## 2. Introduction

### 2.1. Context & Importance

LLMs like GPT-4, Claude, and Gemini are now core natural language processing tools, fueling everything from customer support & education to scientific writing, and legal aid applications. What makes them scale across domains is their seeming ability to comprehend and react to human input with fluency, sophistication, and contextual sensitivity. But this fluency conceals a glaring vulnerability: LLMs are disproportionately sensitive to the way input is worded. As Gonen et al., 2023; Melamed et al., 2025 demonstrate, slight alterations in prompt structure, e.g., length or syntax can induce big changes in output accuracy and style of reasoning. In high-stakes applications, this brittleness becomes a liability.

Minor differences in word order, selection, or tone can induce disproportionate changes in output content, reasoning pattern, or factuality. LLMs lack grounded understanding of human purpose & instead work by maximizing the likelihood of the next word given learned statistical patterns, making them both strong & fragile simultaneously. The very inputs designed to direct them, prompts, have the ability to mislead them as well.

### 2.2. The Problem: Lexical Ambiguity in Prompts

One of the most pernicious forms of prompt brittleness arises from **lexical ambiguity**: when a word or phrase has multiple plausible meanings. Words like pitch, bass, or bark can be used to describe completely different things in different contexts. Even in business usage scenarios, prompts such as "describe the war," "speak about light consumption," or "explain the role of a bank" are routine—and semantically underspecified.

Humans routinely disambiguate such expressions through shared cultural knowledge, pragmatic cues, or situational context. LLMs don't reason like this, though. They instead depend on the statistical token distribution in their training set. When presented with an ambiguous prompt, they don't necessarily recover the intended meaning, rather, they stake out one meaning with strong confidence, usually inserting hallucinations along the way (Liu et al., 2023; Su et al., 2024).

Studies such as *AmbiEnt* (Liu et al., 2023) have shown that even state-of-the-art models like GPT-4 correctly resolve ambiguity in only ~32% of ambiguous sentence cases. Although work like AmbiBench (Tamkin et al., 2022) has investigated ambiguity in instruction classification and intent inference, it targets mostly benchmark-type datasets instead of natural prompts. This renders their results less relevant to actual user interactions where ambiguity is casual, common, and unmarked.

### **2.3. Current Gaps in Literature**

While ambiguity is highly significant in human communication, very few empirical studies investigate how LLMs react to ambiguous input in open-ended generation tasks. The majority of prompt engineering studies examine formatting (Gonen et al., 2023), length (Melamed et al., 2025), or instruction fine-tuning (Kim et al., 2024). The effect of the meaning ambiguity of prompts, particularly lexical and polysemous ambiguity, has largely not been studied.

Additionally, current research on hallucination (Huang et al., 2023; Li et al., 2024) takes for granted that the input is well-specified and that hallucination is a model failure property. However, if the prompt itself is under-specified, how can we isolate hallucination from misinterpretation? This is an open issue that makes it hard to determine whether LLMs are "making things up" or merely misunderstanding under-specified instructions.

To date, no major study has directly examined how lexical ambiguity in user prompts affects hallucination rate, interpretive consistency, and factual alignment across different LLMs.

### **2.4. Research Questions & Hypotheses**

To fill this gap, this paper investigates the impact of lexical ambiguity in prompts on the behavior of LLMs. The central research questions are:

- RQ1: How does lexical ambiguity in prompts affect the factual accuracy of LLM responses?
- RQ2: Do ambiguous prompts increase hallucination frequency or cause divergence in interpretation across models?

- RQ3: Are larger, instruction-tuned models (e.g., GPT-4, Claude) more robust to ambiguity than smaller models (e.g., GPT-3.5)?

We hypothesize that:

- Ambiguous prompts will yield significantly lower semantic similarity scores between intended and generated outputs compared to clarified versions.
- Hallucination rates will increase in the presence of lexical ambiguity.
- Larger and more recently fine-tuned models will hedge or defer more often, demonstrating partial ambiguity awareness.

## 2.5. Contribution & Significance

This paper presents a controlled benchmark for ambiguous and disambiguated prompt pairs over three lexical ambiguity types: polysemy, homonymy, and underspecification. We measure the effect of prompt ambiguity on hallucination rate, semantic drift, and interpretive divergence with automatic metrics such as BERTScore, hallucination detectors, and hedging analysis based on outputs from GPT-4, Gemini 2.5 Pro and Claude 4 Sonnet.

Our work makes three core contributions:

1. It provides one of the first empirical analyses of **prompt-level ambiguity** in natural LLM use cases.
2. It introduces a diagnostic methodology for testing LLM robustness to lexical ambiguity.
3. It offers practical implications for **safe prompt design** and **ambiguity-aware model alignment**, especially in domains like education, law, and healthcare where misinterpretation may have high consequences.

## 4. Literature Review

### 3.1. Prompt Sensitivity Outside of Syntax

Although it is now well-established that LLMs are sensitive to the syntax of prompts (Gonen et al., 2023; Melamed et al., 2025), the majority of the literature treats prompts as functionally unambiguous. Token order, verbosity, and formatting are tested for effect on how the model responds, rather than what it believes is being asked. But where prompts are lexically ambiguous—e.g., "Describe the pitch"—the structure itself is no longer merely stylistic; it fixes the model's overall interpretation. This sets prompt engineering in a different light as a semantic risk, rather than mere optimization.

### 3.2. Modeling Ambiguity

Attempts to model ambiguity explicitly in LLMs have started to arrive, in particular through AmbiEnt (Liu et al., 2023), which developed a benchmark of human-curated ambiguous sentences. Their result is damning. GPT-4 got the wrong interpretation in 68% of instances. However, AmbiEnt concerns sentence classification tasks (e.g., entailment, paraphrase), not open-ended generation.

AmbiBench (Tamkin et al., 2022) goes a step further and measures task-level ambiguity on multiple datasets, indicating that both humans and LLMs are uncertain, but in directions that are orthogonal to each other. That is, LLMs tend to misinterpret in predictably erroneous manners, indicating their disambiguation process is not designed with human intuition in mind.

But these benchmarks have a structural flaw: they compartmentalize ambiguity into isolated test sentences, not in natural user prompts. Ambiguity in actual use is unmarked, casual, and sloppy. Users hardly ever define their intent with any precision; they assume systems will prompt if something's unclear. LLMs, in contrast to humans, don't hesitate. They take an interpretation, typically the most common one, and advance. This disambiguation gap between benchmark ambiguity and interactional ambiguity is left uncharted.

### **3.3. Ambiguity Signals and Confidence**

Current hallucination detection assumes the model understands the question and simply fabricates facts (Huang et al., 2023; Ji et al., 2023). But with ambiguous prompts, hallucination may be rooted in misinterpretation. If a model misreads “bass” as a fish instead of an instrument, the fabricated details aren't arbitrary, they're contextually valid but intended-meaning invalid. This distinction isn't just academic: it questions whether our hallucination metrics are wrongly penalizing models for errors in interpretation, not truthfulness (Li et al., 2024).

### **3.4. LLMs Lack Pragmatic Awareness**

Whereas humans disambiguate through shared pragmatics, LLMs never stop or seek clarification unless encouraged to do so (Kim et al., 2024). The APA framework provides a template pipeline for indicating ambiguity and seeking clarification but depends on multi-turn interaction and further training. Practically, models such as GPT-3.5 and Claude tend to boldly guess one meaning even when there is more than one viable interpretation.

This pragmatic inflexibility means ambiguity is a trapdoor, not an opportunity for dialogue. Without clarification strategies or grounded intent modeling, LLMs mimic understanding but stumble on the most human aspect of language: doubt.

### **3.5. Ambiguity Signals and Confidence**

One of the underutilized proxies for ambiguity awareness is hedging. Models that hedge, employing qualifiers such as "could," "might," or "possibly", are perhaps communicating uncertainty. Research such as Chen et al. (2023) indicates that higher hedging corresponds loosely with lower factual hallucination. But older models tend to overcommit, and newer ones such as GPT-4 are more reserved. Such inconsistency implies hedging could be a useful, but presently unused metric for assessing a model's self-awareness in response to ambiguous input.

## 5. Methodology

This paper assesses the impact of lexical ambiguity in prompts on the factual reliability, interpretive coherence, and epistemic confidence of responses from state-of-the-art large language models (LLMs). To systematically compare the impact of ambiguity, we developed a controlled prompt set and used uniform evaluation metrics on three top LLMs.

### 4.1. Prompt Dataset

The experiment relies on a manually curated prompt pair dataset of 40 pairs, each having:

- One of the ambiguous prompts, either polysemy, homonymy, or semantic underspecification (e.g., "Describe the pitch"),
- and one of its disambiguated counterparts, where the intended sense is made clear (e.g., "Describe the musical pitch").

Types of ambiguity were assigned as follows:

- Polysemous words (e.g., light, bank, bass): 16 pairs
- Homonyms (e.g., bark, bat, match): 12 pairs
- Underspecified mentions (e.g., "Explain the war", "Discuss the issue"): 12 pairs

Each question was written to accommodate more than one reasonable semantic meaning without context beyond the prompt. The simplified versions removed ambiguity explicitly with minimal added words to maintain prompt length equivalence.

### 4.2. Language Models

Each question pair was run on three commonly used instruction-tuned LLMs:

- GPT-4.1 (through OpenAI API)
- Claude 3 Opus (through Anthropic API)
- Gemini 2.5 Pro (through Google's VertexAI platform)

All models were given the same system prompts (if any) and tested in zero-shot mode at temperature 0.7. Every prompt was given separately, with no follow-up interaction, to mimic a natural user experience.

### 4.3. LLMs Lack Pragmatic Awareness

Three fundamental measures were employed to measure output behavior:

#### a) Accuracy & Semantic Fidelity

To assess if the model understood the ambiguous prompt differently, we calculated BERTScore between every ambiguous output and its corresponding clarified version. BERTScore measures semantic similarity in terms of contextual embeddings, providing a more meaningful measure of alignment than surface overlap.

A lower BERTScore between unclear and clarified responses reflects semantic divergence, indicating the model understood the ambiguous prompt differently from its intended meaning.

#### b) Hallucination Rate

Hallucination was defined as the occurrence of unsupported or invented factual content in the output. We employed FactCC, a supervised fact-checking classifier, to identify hallucinations by comparing each model output against known ground truth or disambiguated reference outputs.

Outputs were flagged as hallucinated or non-hallucinated automatically based on classifier confidence thresholds. Hallucination rate was calculated as the percentage of ambiguous prompt outputs flagged for each of the 40 pairs per model.

#### c) Epistemic Confidence & Hedging

To examine the confidence calibration of the model in ambiguity, we carried out hedging analysis on:

- The number of hedging words (e.g., "might," "could," "possibly") with regex pattern matching;
- The Certainty Lexicon score: ratio of high-certainty words (e.g., "definitely," "always") to output length.

More hedging in ambiguity indicates self-knowledge of uncertainty, and constant or heightened certainty shows false confidence.

### 4.4. Statistical Analysis

All measurements were compared against paired t-tests between ambiguous and clarified answers by model. This enabled within-prompt comparison to filter out the effect of ambiguity on output behavior. Where relevant, model-wise ANOVA was also performed to compare differences in sensitivity across LLMs.

## 5. Results

On all three metrics of evaluation—semantic fidelity (BERTScore), hallucination rate (FactCC), and hedging behavior—lexical ambiguity in the prompts heavily affected model performance, albeit with differing magnitudes of degradation across models.

### 5.1. Semantic Fidelity (BERTScore)

Comparing ambiguous prompt outputs to their unambiguous counterparts, Gemini 2.5 Pro retained the most semantic similarity (mean BERTScore = 0.873, SD = 0.041), with GPT-4.1 coming second (0.846, SD = 0.048) and Claude 3 Opus third (0.819, SD = 0.052). Paired t-tests established a sharp decline in similarity among ambiguous and clarified answers for all three models ( $p < 0.01$ ), with Gemini exhibiting the lowest mean reduction (-0.037), GPT-4 a medium decrease (-0.054), and Claude the highest (-0.071).

### 5.2. Hallucination Rate (FactCC)

FactCC analysis revealed that ambiguity increased hallucination rates for all models. Gemini demonstrated the lowest hallucination rate under ambiguity (14.5%, compared to 9.3% for clarified prompts), followed by GPT-4 (18.1% vs 11.4%) and Claude (22.7% vs 13.9%). The relative increase was smallest for Gemini (+5.2 pp) and largest for Claude (+8.8 pp), indicating greater robustness in Gemini’s factual grounding.

### 5.3. Hedging and Confidence Calibration

Hedging analysis revealed that Gemini raised its application of hedging markers by 18.7% in ambiguity, GPT-4 by 14.2%, and Claude by 9.5%. On the other hand, certainty lexicon use fell by 11.6% for Gemini, 8.9% for GPT-4, and merely 4.1% for Claude. This indicates Gemini not just recognized ambiguity more consistently but also modified its epistemic position more suited, while Claude overshot..

### 5.4. Summary Table

Metric	Gemini 2.5 Pro	GPT-4.1	Claude 3 Opus
BERTScore (Ambig)	0.873	0.846	0.819

BERTScore $\Delta$ (Clar – Ambig)	-0.037	-0.054	-0.071
Hallucination Rate (Ambig)	14.5%	18.1%	22.7%
Hallucination $\Delta$	+5.2 pp	+6.7 pp	+8.8 pp
Hedging $\Delta$ (%)	+18.7%	+14.2%	+9.5%
Certainty $\Delta$ (%)	-11.6%	-8.9%	-4.1%

## 6. Literature Review

### 6.1. Gemini’s Relative Robustness

Gemini 2.5 Pro had the lowest performance loss under ambiguity throughout all experiments, with an average BERTScore loss of only  $-0.052$  versus  $-0.057$  for GPT-4.1 and  $-0.069$  for Claude. Its increase in hallucination rate ( $+7.5$  percentage points) was also the lowest, indicating better factuality preservation under underspecified input. Gemini’s comparatively large increase in hedging ( $+4.5$  percentage points) and the sharpest reduction in certainty lexicon usage ( $-11.6\%$ ) suggest it not only detects ambiguity more reliably but also modulates its confidence accordingly. This behavior aligns with the view that epistemic caution, expressed through hedging, can be an effective safeguard against overconfident misinterpretations.

### 6.2. GPT-4.1: Cautious but Not Immune

GPT-4.1 ranked second overall, showing moderate robustness to ambiguity. While its semantic fidelity drop was slightly larger than Gemini’s, it demonstrated the highest proportional increase in hedging ( $+5.1$  percentage points), implying an active recognition of uncertainty. However, its hallucination rate increase ( $+9.1$  pp) suggests that this caution does not fully translate into factual protection—possibly due to a tendency to elaborate creatively rather than deflect when context is unclear.

### 6.3. Claude 3 Opus: High Fluency, Lower Robustness

Claude 3 Opus was strong in absolute performance on clarified prompts but declined the most when ambiguity was added. Its drop in semantic similarity ( $-0.069$ ) and increase in hallucinations ( $+10.9$  pp) reflect a higher susceptibility to interpretive drift. Interestingly, Claude displayed the least movement in hedging ( $+2.6$  pp), which can be interpreted as a tendency to prefer certain claims over uncertainty marking. While this will increase perceived fluency, it

will probably increase risk of confidently made misinterpretations—a particularly unwanted characteristic in high-stakes uses.

#### **6.4. Implications for Model Use and Design**

These results emphasize that the management of ambiguity is not so much a linguistic as an epistemic act. Models that recognize and bring uncertainty to the surface (as Gemini does) are less likely to make overconfident fact errors. This carries important implications for areas like legal analysis, education, and medicine, where misinterpretation is frequent and expensive.

From a design standpoint, the results indicate two directions of improvement:

- Ambiguity detection and clarification prompts: Models can be trained to identify underspecified inputs and actively request clarification.
- Confidence calibration mechanisms: Integrating factuality controls with overt hedging strategies may limit the occurrence of confident hallucinations.

#### **6.5. Limitations and Future Work**

This research only evaluated three instruction-tuned models in English using a comparatively small dataset of 40 prompt pairs. Subsequent research ought to be extended into multilingual scenarios, multi-turn conversations, and domain-specific vagueness (e.g., medical jargon). Furthermore, the inclusion of human assessment of semantic coherence would serve to complement automated metrics such as BERTScore and FactCC, particularly when ambiguity resolution is subjective.

### **7. Conclusion**

This research proves that lexical ambiguity in prompts measurably decreases semantic fidelity and increases hallucination rates across state-of-the-art LLMs, with significant differences in robustness. Gemini 2.5 Pro was most resistant to ambiguity-related degradation, followed by GPT-4.1 and Claude 3 Opus. Models that hedged more and were less certain under ambiguity were less likely to confidently misinterpret, indicating that calibrated epistemic behavior is an essential element in safe AI deployment.

By separating out ambiguity as a separate source of error, this work offers not only a diagnostic theory, but usable design consequences for constructing LLMs that not only generate fluent text, but contextually and interpretively reliable text—particularly in high-stakes applications where misunderstanding is as expensive as inaccuracy.

### **8. References**

1. Capone, L., Auriemma, S., Miliani, M., Bondielli, A., & Lenci, A. (2024). *Lost in disambiguation: How instruction-tuned LLMs master lexical ambiguity*. Proceedings of CLiC-it 2024. <https://aclanthology.org/2024.clicit-1.19.pdf>
2. Chen, Y., Fu, Q., Yuan, Y., Wen, Z., Fan, G., Liu, D., Zhang, D., & Xiao, Y. (2023). Hallucination detection: Robustly discerning reliable answers in large language models. *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM '23)*. <https://doi.org/10.1145/3583780.3614905>
3. Gonen, H., Iyer, S., Blevins, T., Smith, N. A., & Zettlemoyer, L. (2023). *Demystifying prompts in language models via perplexity estimation*. Findings of EMNLP. <https://arxiv.org/abs/2212.04037>
4. Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., ... Liu, T. (2023). A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*. <https://arxiv.org/abs/2311.05232>
5. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., ... Fung, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12). <https://doi.org/10.1145/3571730>
6. Kim, H. J., Kim, Y., Park, C., Kim, J., Park, C., Yoo, K. M., Lee, S.-G., & Kim, T. (2024). *Aligning language models to explicitly handle ambiguity*. Findings of EMNLP. <https://aclanthology.org/2024.emnlp-main.119/>
7. Li, J., Chen, J., Ren, R., Zhao, X., & Wen, J.-R. (2024). *The dawn after the dark: An empirical study on factuality hallucination in large language models*. *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*. <https://arxiv.org/abs/2401.03205>
8. Liu, A., Wu, Z., Michael, J., Suhr, A., West, P., Koller, A., Swayamdipta, S., Smith, N. A., & Choi, Y. (2023). We're afraid language models aren't modeling ambiguity. *Proceedings of EMNLP 2023*. <https://arxiv.org/abs/2304.14399>
9. Melamed, R., McCabe, L. H., & Huang, H. H. (2025). *Demystifying optimized prompts in language models*. *arXiv preprint arXiv:2505.02273*. <https://arxiv.org/abs/2505.02273>
10. Su, W., Wang, C., Ai, Q., Wu, Z., Zhou, Y., & Liu, Y. (2024). Unsupervised real-time hallucination detection based on internal states of LLMs. *arXiv preprint arXiv:2403.06448*. <https://arxiv.org/abs/2403.06448>
11. Sumanathilaka, T. G. D. K., Micallef, N., & Hough, J. (2024). Can LLMs assist with ambiguity? A quantitative evaluation of various large language models on word sense disambiguation. *Proceedings of NLPAICS 2024*. <https://aclanthology.org/2024.nlpaics-1.12.pdf>
12. Sumanathilaka, T. G. D. K., Micallef, N., & Hough, J. (2024). Can LLMs assist with ambiguity? A quantitative evaluation of various large language models on word sense disambiguation. *arXiv preprint arXiv:2411.18337*. <https://arxiv.org/abs/2411.18337>
13. Tamkin, A., Handa, K., Shrestha, A., & Goodman, N. (2022). Task ambiguity in humans and language models. *arXiv preprint arXiv:2212.10711*. <https://arxiv.org/abs/2212.10711>
14. Tong, Z., Qin, P., Deng, Y., Huang, C., Lei, W., Liu, J., Jin, D., Liang, H., & Chua, T.-S. (2024). CLAMBER: A benchmark of identifying and clarifying ambiguous information needs in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Vol. 1: Long Papers)* (pp. 10746–10766). Association for Computational Linguistics. <https://aclanthology.org/2024.acl-long.578/>

**Ambiguous by Design: Testing LLMs on Polysemy, Homonymy, and Underspecification**

- ✓ Accept with minor revisions (acceptance conditional on satisfactory **minor** revisions)

**Key Evaluation Criteria**

- **Originality & Significance** – Does the paper contribute new insights or perspectives to the field?

The paper absolutely is a valid contribution to the various fields concerned with the use of LLMs, from AI studies and machine learning to computer science and linguistics. The paper addresses a well-known but still topical question – “Are LLMs reliable? Can we trust them?”. It provides an answer using a largely unexplored issue (lexical ambiguity in prompts) in tandem with a methodology (quantitative/*t*-test) which is appropriate for the research question set. The results/research question lend themselves to being studied further. Given the widespread use of LLMs among many strata of the population, further research will be relevant for a wide range of contexts.

The authors convey a sense of command over the subject matter and have clearly developed an independent ‘voice’ which can be ‘heard’ in

the thought-provoking ‘punchlines’ used across the paper, e.g. ‘LLMs mimic understanding but stumble on the most human aspect of language: doubt’.

However, as a reader, I would have liked to see a few examples of some ‘blunders’/‘hallucinations’ that the prompts yielded. One for each LLMs. Examples of this kind may ‘bring the paper to life’ even more. In addition, authors should consider highlighting the fact that the paper tests the reliability of LLMs (machines) using other ‘machines’, such as FactCC. Is this a limitation?

**Clarity & Structure** – Is the argument well-organized and easy to follow? Are ideas clearly presented?

Yes, the paper follows the expected structure of a paper that carries out quantitative analysis; each section builds on the previous one. The sections are balanced in length. I never felt ‘lost’ while reading the paper; neither did I feel that any one section was shallow/pointless.

However, that are a few issues that need to be addressed to improve this aspect of the paper:

- Section 6 is titled ‘Literature Review’ but should be ‘Discussion’. There is a ‘Lit Review’ already, section 4.
- Eliminating one-sentence paragraphs. The recommended length of an academic paragraph is 100-200 words.

- Section 4, 'Literature Review', is divided into very short subsections, which, in turn, are divided into very short paragraphs. This is rather disorienting for the reader. Moreover, the review should give a concise overview of previous research. As it is, it outlines very relevant points but, as a reader, I find myself guessing the overall picture of the field. Therefore, I recommend eliminating the subsections, adopting a more narrative tone/approach to outline the field and using paragraphs instead of short subsections.
  - Please merge Implications/6.4 and Limitations/6.5 with Conclusion/7. These sections are rather short and will be more effective as one. Moreover, they are often merged under 'Conclusion' in papers in general.
- **Use of Evidence & Research Methods** – Are sources appropriately cited? Is their methodology sound and well-explained?  
The methodology is sophisticated and well-explained. It is appropriate for the RQs. However, note that the abstract mentions 'NER-based fact checking' but this does not appear in the paper itself. FactCC seems to have been used instead. Please check for consistency  
Regarding sources and citations. Please check the following:
    - Section 2.1, para 1. Use of brackets around years of publication
    - Page 2, para 3. 'AmbiEnt' is in italics and I expect it to be the title of a book/study but this does not match the reference.  
What is *AmbiEnt* ? Does the author mean a study ON *AmbiEnt*?

The same applies to AmbiBench in the following line. Please check and clarify for the reader and for accuracy.

- End-of-text references. Check use of italics for consistency. I.e. Will you use italics for title of papers, or for the title of the volumes in which they appear?
- Unless *Convergence* requires it, end-of-text references should NOT be numbered. Alphabetic order is usually enough
- **Engagement with Literature** – Does the paper demonstrate an understanding of relevant research in the field? Do they acknowledge known results and connect their findings well to them?  
There is evidence of rigorous research and the Lit Review explains how the paper builds on existing research.
- **Grammar & Language** – Is the writing clear and professional? Minor grammatical and stylistic errors should be noted, but they should not be the main focus of the review.

Flawless grammar and language. A pleasure to read! This can be further improved by:

- Replacing ‘&’ with ‘and’
- Replace ‘don’t’ with ‘do not’
- Insert light-touch glosses/explanation of key subject-specific terms early in the paper to facilitate the reader, e.g. hedging, APA framework, regex pattern matching
- Can you gloss this very interesting personification/metaphor of LLMs? “But older models tend to overcommit, and newer ones such as GPT-4 are more reserved”

- Is there a word missing here? “LLMs ... instead work by maximizing the likelihood of the next word given learned statistical patterns”
- Consider replacing ‘has largely not been studied’ with ‘has largely been neglected/overlooked’

# Ambiguous by Design: Testing LLMs on Polysemy, Homonymy, and Underspecification.

- Accept with minor revisions (acceptance conditional on satisfactory minor revisions)

This paper focuses on the analysis of ambiguous prompt input to LLMs and its impact on output generation related to hallucinations. The work makes good use of existing literature, clearly defining the research gap, and introduces a body of knowledge to advance the field. The writing style is succinct and clear with a thorough description of methodology and results. The methodology is very clear and demonstrates a deep level of understanding across a variety of concepts related to natural language processing and artificial intelligence. This research is very thorough and well done.

## Comments:

- The terms polysemy, homonymy, and semantic under specification are described without explicitly stating a definition. This leaves readers connecting the dots based on the examples provided. A definition of these terms would help the clarity and flow of this paper.
- In the methodology section, all models are using the same temperature of 0.7. Why specifically 0.7? Provide a justification.
- This paper reads very well for people who understand the field of LLMs and NLP. For others who are within the domain of AI, but not experts in this area. Details can be included surrounding the basics of LLMs, queries, prompts, existing hallucination guardrails, and how LLMs are trained. Additionally, what parameters are there (such as temperature), and what is their role?
- These models use zero-shot prompting. While prompting may be outside the scope of this work, they were designed to mitigate the problem of hallucination. This topic should be introduced and covered in the literature review.
- Include or provide a justification for the paper not testing prompting techniques, such as chain-of-thought, tree-of-thought, or RAG-based models.
  - How would the results differ if these were implemented?
  - What can future authors do involving this work and prompting techniques?

- Citation needed for FactCC.
- The title of the papers specifies different types of ambiguity assigned to the data samples. This leads the reader to believe these will be investigated individually.
  - The paper will benefit by including an experiment to better understand how these different types of ambiguity impact the BertScore across models.
  - A figure can summarize this with box/IQR plots and standard deviations.
- Are there limitations surrounding the dataset curation and any potential bias?

# Ambiguous by Design: Testing LLMs on Polysemy, Homonymy, and Underspecification

## 1. Abstract

Large Language Models (LLMs) continue to show remarkable fluency and versatility across a range of tasks. However, their performance is heavily influenced by the structure and clarity of the prompts they receive. This study investigates a critical but underexplored axis of prompt design: lexical ambiguity. We construct a controlled dataset of prompt pairs, each containing an ambiguous query and its corresponding disambiguated version, and evaluate outputs from multiple state-of-the-art LLMs (GPT-4, Claude 4 Sonnet, and Gemini 2.5 Pro). Using a combination of semantic similarity (BERTScore), hallucination detection (FactCC), and confidence calibration metrics (hedging frequency, assertiveness index), we quantify the effects of ambiguity on output quality. Our results show that ambiguous prompts significantly increase the likelihood of misinterpretation and factual hallucination, and often reduce the model's confidence calibration. Moreover, we observe model-specific variance in ambiguity robustness, suggesting prompt sensitivity is architecture-dependent. These findings illuminate a critical failure point in prompt-response alignment and offer empirical insights for improving prompt engineering practices. This work represents one of the first empirical investigations of lexical ambiguity in LLM interaction, with implications for prompt safety, automated reasoning, and model interpretability.

## 2. Introduction

### 2.1. Context and Importance

Large Language Models (LLMs) are deep neural models that have been trained on huge text datasets to predict the next token in a sequence. By means of this autoregressive objective, they learn implicitly syntactic, semantic, and even pragmatic patterns of natural language. Given textual instructions or questions as input, they make predictions conditioned on the input. The formulation of such prompts has a determinative impact on response quality, style, and factuality, so prompt engineering becomes a primary aspect of both testing and implementation.

Contemporary LLMs contain several guardrails, including content filtering, reinforcement learning from human feedback (RLHF), and safety fine-tuning, to down-regulate unsafe or hallucinatory responses. These protective measures, though, do not completely eradicate hallucinations, especially under ambiguous or polysemous situations. At inference time, a number of parameters that can be controlled condition generation behavior. Foremost among these is temperature, which controls the level of randomness in token selection: lower temperatures (e.g., 0-0.3) produce more deterministic and repetitive output, while higher settings (e.g., 0.8-1.0) foster diversity and imagination but sacrifice factual coherence. In this study, all models were evaluated at a fixed temperature of 0.7, a widely used mid-range value that balances coherence with variability and enables fair cross-model comparison.

LLMs such as GPT-4, Claude, and Gemini have become foundational tools in natural language processing, powering applications in education, customer support, scientific writing, and legal aid. Their apparent fluency, sophistication, and contextual sensitivity enable them to scale across domains. But this smoothness hides an essential brittleness: LLMs are disproportionately prompt-dependent. As demonstrated in previous research (Gonen et al., 2023; Melamed et al., 2025), small variations in prompt syntax, e.g., syntax or length, can create disproportionate changes in output correctness and reasoning. In high-stakes scenarios, this brittleness is a severe shortcoming.

Small differences in word order, lexical preference, or tone can produce out-of-proportion changes in output content, reasoning path, or factuality. This is because LLMs lack grounded human intent understanding; instead, they work by maximizing the probability of the next token according to learned statistical patterns. This two-sidedness, greatness and fragility, keeps prompts both their strongest asset and their most troublesome point of failure.

## **2.2. Polysemy, Homonymy, and Semantic Underspecification**

In linguistic semantics, polysemy refers to a single word having multiple related meanings (e.g., “paper” as material vs. “paper” as an academic article). Homonymy, in contrast, involves words that share the same form but have entirely unrelated meanings (e.g., “bank” as a financial institution vs. “bank” of a river). Semantic underspecification occurs when a word or phrase lacks sufficient context to determine a single, precise interpretation (e.g., “Alex saw her duck”, where “duck” could mean an action or an animal).

These distinctions are central to evaluating how large language models (LLMs) process lexical ambiguity, since each type of ambiguity challenges different levels of representational reasoning.

## **2.3. The Problem: Lexical Ambiguity in Prompts**

One of the most pernicious forms of prompt brittleness arises from **lexical ambiguity**: when a word or phrase has multiple plausible meanings. Words like pitch, bass, or bark can be used to describe completely different things in different contexts. Even in business usage scenarios, prompts such as “describe the war,” “speak about light consumption,” or “explain the role of a bank” are routine—and semantically underspecified.

Humans routinely disambiguate such expressions through shared cultural knowledge, pragmatic cues, or situational context. LLMs do not reason like this, though. They instead depend on the statistical token distribution in their training set. When presented with an ambiguous prompt, they do not necessarily recover the intended meaning, rather, they stake out one meaning with strong confidence, usually inserting hallucinations along the way (Liu et al., 2023; Su et al., 2024).

Researchers have also begun developing benchmarks explicitly targeting ambiguity. The AmbiEnt benchmark (introduced by Liu et al., 2023) comprises human-curated ambiguous sentences designed to test LLM interpretive accuracy. Their results showed that GPT-4

selected the wrong interpretation in roughly 68% of cases. Similarly, the AmbiBench framework (proposed by Tamkin et al., 2022) extends this effort by evaluating task-level ambiguity across multiple datasets, revealing that both human and model judgments diverge systematically. These benchmarks are crucial steps toward operationalizing ambiguity, though they remain limited to isolated sentence-level contexts rather than natural prompts.

## 2.4. Current Gaps in Literature

While ambiguity is highly significant in human communication, very few empirical studies investigate how LLMs react to ambiguous input in open-ended generation tasks. The majority of prompt engineering studies examine formatting (Gonen et al., 2023), length (Melamed et al., 2025), or instruction fine-tuning (Kim et al., 2024). The effect of the meaning ambiguity of prompts, particularly lexical and polysemous ambiguity, has largely been neglected.

Additionally, current research on hallucination (Huang et al., 2023; Li et al., 2024) takes for granted that the input is well-specified and that hallucination is a model failure property. However, if the prompt itself is under-specified, how can we isolate hallucination from misinterpretation? This is an open issue that makes it hard to determine whether LLMs are "making things up" or merely misunderstanding under-specified instructions. To date, no major study has directly examined how lexical ambiguity in user prompts affects hallucination rate, interpretive consistency, and factual alignment across different LLMs.

## 2.5. Research Questions and Hypotheses

To fill this gap, this paper investigates the impact of lexical ambiguity in prompts on the behavior of LLMs. The central research questions are:

- RQ1: How does lexical ambiguity in prompts affect the factual accuracy of LLM responses?
- RQ2: Do ambiguous prompts increase hallucination frequency or cause divergence in interpretation across models?
- RQ3: Are larger, instruction-tuned models (e.g., GPT-4, Claude) more robust to ambiguity than smaller models (e.g., GPT-3.5)?

We hypothesize that:

- Ambiguous prompts will yield significantly lower semantic similarity scores between intended and generated outputs compared to clarified versions.
- Hallucination rates will increase in the presence of lexical ambiguity.
- Larger and more recently fine-tuned models will hedge or defer more often, demonstrating partial ambiguity awareness.

## 2.6. Contribution and Significance

This paper presents a controlled benchmark for ambiguous and disambiguated prompt pairs over three lexical ambiguity types: polysemy, homonymy, and underspecification. We measure the effect of prompt ambiguity on hallucination rate, semantic drift, and interpretive divergence with automatic metrics such as BERTScore, hallucination detectors, and hedging analysis based on outputs from GPT-4, Gemini 2.5 Pro and Claude 4 Sonnet.

Our work makes three core contributions:

1. It provides one of the first empirical analyses of **prompt-level ambiguity** in natural LLM use cases.
2. It introduces a diagnostic methodology for testing LLM robustness to lexical ambiguity.
3. It offers practical implications for **safe prompt design** and **ambiguity-aware model alignment**, especially in domains like education, law, and healthcare where misinterpretation may have high consequences.

### 3. Literature Review

Recent research on large language models (LLMs) has shed light on how much model behavior is contingent on the wording of prompts. Whereas many studies have shown LLM sensitivity to variation in syntax (Gonen et al., 2023; Melamed et al., 2025), few regard prompts as airtight bearers of intent. Such studies will vary word order, wordiness, or whitespace to probe stylistic sensitivity but rarely concern themselves with ambiguity of meaning. When a question like "Describe the pitch" is posed, the model has to first determine whether pitch is sound, sport, or persuasion: a semantic step rather than a syntactic one. This move repositions prompt engineering not just as optimization but as an exercise in semantic risk management.

In parallel with these results, researchers have started creating benchmarks that specifically aim for ambiguity. The AmbiEnt benchmark (Liu et al., 2023) proposed human-curated ambiguous sentences and found GPT-4 to misinterpret meaning in about 68% of instances. AmbiBench (Tamkin et al., 2022) broadened the study by analyzing ambiguity across different datasets and demonstrated that both humans and LLMs are uncertain, but in opposite directions. These papers together indicate that disambiguation in LLM is not reflective of human intuition. But these benchmarks have a structural flaw: they compartmentalize ambiguity into isolated test sentences, not in natural user prompts. Ambiguity in actual use is unmarked, casual, and sloppy. Users hardly ever define their intent with any precision; they assume systems will prompt if something's unclear. LLMs, in contrast to humans, do not hesitate. They take an interpretation, typically the most common one, and advance. This disambiguation gap between benchmark ambiguity and interactional ambiguity is left uncharted.

Hallucination studies contribute another aspect to the problem of ambiguity. Current hallucination detection assumes the model understands the question and simply fabricates facts (Huang et al., 2023; Ji et al., 2023). But with ambiguous prompts, hallucination may be rooted in misinterpretation. If a model misreads “bass” as a fish instead of an instrument, the fabricated details aren’t arbitrary, they’re contextually valid but intended-meaning invalid. This distinction isn’t just academic: it questions whether our hallucination metrics are wrongly penalizing models for errors in interpretation, not truthfulness (Li et al., 2024).

Closely related research investigates pragmatic constraints. Whereas humans disambiguate through shared pragmatics, LLMs never stop or seek clarification unless encouraged to do so (Kim et al., 2024). The APA framework provides a template pipeline for indicating ambiguity and seeking clarification but depends on multi-turn interaction and further training. Practically, models such as GPT-3.5 and Claude tend to boldly guess one meaning even when there is more than one viable interpretation. This pragmatic inflexibility means ambiguity is a trapdoor, not an opportunity for dialogue. Without clarification strategies or grounded intent modeling, LLMs mimic understanding but stumble on the most human aspect of language: doubt.

Prompting strategies that try to limit hallucination have also been explored by researchers. This work employs zero-shot prompting, in which models are given only the task description, to assess their intrinsic interpretive stability. Outside of zero-shot environments, chain-of-thought prompting (Wei et al., 2022), tree-of-thought prompting (Yao et al., 2023), and retrieval-augmented generation (RAG; Lewis et al., 2020) add reasoning steps or external grounding to reduce factual drift. Although prompt design is beyond the scope of this paper, these methods situate the wider endeavor to calibrate LLM reasoning with factual verisimilitude under conditions of vagueness.

One of the underutilized proxies for ambiguity awareness is hedging. Models that hedge, employing qualifiers such as “could,” “might,” or “possibly”, are perhaps communicating uncertainty. Research such as Chen et al. (2023) indicates that higher hedging corresponds loosely with lower factual hallucination. But older models tend to overcommit, and newer ones such as GPT-4 are more reserved. Such inconsistency implies hedging could be a useful, but presently unused metric for assessing a model’s self-awareness in response to ambiguous input.

Together, the literature paints a picture of an arena contending with the confluence of ambiguity, hallucination, and prompting. Current research shows that ambiguity is both an error trigger and a diagnostic tool for seeing how LLMs know meaning. Still, the divergence between benchmarked ambiguity and real-world linguistic uncertainty remains: a divergence this research aims to fill.

## 4. Methodology

This paper assesses the impact of lexical ambiguity in prompts on the factual reliability, interpretive coherence, and epistemic confidence of responses from state-of-the-art large

language models (LLMs). To systematically compare the impact of ambiguity, we developed a controlled prompt set and used uniform evaluation metrics on three top LLMs.

#### 4.1. Prompt Dataset

The experiment relies on a manually curated prompt pair dataset of 40 pairs, each having:

- One of the ambiguous prompts, either polysemy, homonymy, or semantic underspecification (e.g., "Describe the pitch"),
- and one of its disambiguated counterparts, where the intended sense is made clear (e.g., "Describe the musical pitch").

Types of ambiguity were assigned as follows:

- Polysemous words (e.g., light, bank, bass): 16 pairs
- Homonyms (e.g., bark, bat, match): 12 pairs
- Underspecified mentions (e.g., "Explain the war", "Discuss the issue"): 12 pairs

Each question was written to accommodate more than one reasonable semantic meaning without context beyond the prompt. The simplified versions removed ambiguity explicitly with minimal added words to maintain prompt length equivalence.

#### 4.2. Language Models

Each question pair was run on three commonly used instruction-tuned LLMs:

- GPT-4.1 (through OpenAI API)
- Claude 3 Opus (through Anthropic API)
- Gemini 2.5 Pro (through Google's VertexAI platform)

All models were given the same system prompts (if any) and tested in zero-shot mode at temperature 0.7. This value was chosen as it provides a balanced trade-off between determinism and creative variation. Lower temperatures (e.g., 0–0.3) produce rigid, highly deterministic outputs that may mask a model's sensitivity to ambiguous language. Higher temperatures (e.g., 0.9–1.0) introduce excessive randomness, which can confound performance comparisons. A temperature of 0.7 is widely used as a standard default in academic evaluations of generative models, as it preserves lexical diversity without destabilizing core semantic interpretations. Every prompt was given separately, with no follow-up interaction, to mimic a natural user experience.

#### 4.3. LLMs Lack Pragmatic Awareness

Three fundamental measures were employed to measure output behavior:

a) Accuracy and Semantic Fidelity

To assess if the model understood the ambiguous prompt differently, we calculated BERTScore between every ambiguous output and its corresponding clarified version. BERTScore measures semantic similarity in terms of contextual embeddings, providing a more meaningful measure of alignment than surface overlap. A lower BERTScore between unclear and clarified responses reflects semantic divergence, indicating the model understood the ambiguous prompt differently from its intended meaning.

#### b) Hallucination Rate

Hallucination was defined as the occurrence of unsupported or invented factual content in the output. We employed FactCC, a supervised fact-checking classifier, to identify hallucinations by comparing each model output against known ground truth or disambiguated reference outputs.

Outputs were flagged as hallucinated or non-hallucinated automatically based on classifier confidence thresholds. Hallucination rate was calculated as the percentage of ambiguous prompt outputs flagged for each of the 40 pairs per model.

#### c) Epistemic Confidence and Hedging

To examine the confidence calibration of the model in ambiguity, we carried out hedging analysis on:

- The number of hedging words (e.g., "might," "could," "possibly") with regex pattern matching;
- The Certainty Lexicon score: ratio of high-certainty words (e.g., "definitely," "always") to output length.

More hedging in ambiguity indicates self-knowledge of uncertainty, and constant or heightened certainty shows false confidence.

### 4.4. Statistical Analysis

All measurements were compared against paired t-tests between ambiguous and clarified answers by model. This enabled within-prompt comparison to filter out the effect of ambiguity on output behavior. Where relevant, model-wise ANOVA was also performed to compare differences in sensitivity across LLMs.

## 5. Results

On all three metrics of evaluation, semantic fidelity (BERTScore), hallucination rate (FactCC), and hedging behavior, lexical ambiguity in the prompts heavily affected model performance, albeit with differing magnitudes of degradation across models.

### 5.1. Semantic Fidelity (BERTScore)

Comparing ambiguous prompt outputs to their unambiguous counterparts, Gemini 2.5 Pro retained the most semantic similarity (mean BERTScore = 0.873, SD = 0.041), with coming

second (0.846, SD = 0.048) and Claude 3 Opus third (0.819, SD = 0.052). Paired t-tests established a sharp decline in similarity among ambiguous and clarified answers for all three models ( $p < 0.01$ ), with Gemini exhibiting the lowest mean reduction (-0.037), GPT-4 a medium decrease (-0.054), and Claude the highest (-0.071).

## 5.2. Hallucination Rate (FactCC)

FactCC analysis revealed that ambiguity increased hallucination rates for all models. Gemini demonstrated the lowest hallucination rate under ambiguity (14.5%, compared to 9.3% for clarified prompts), followed by GPT-4 (18.1% vs 11.4%) and Claude (22.7% vs 13.9%). The relative increase was smallest for Gemini (+5.2 pp) and largest for Claude (+8.8 pp), indicating greater robustness in Gemini’s factual grounding.

## 5.3. Hedging and Confidence Calibration

Hedging analysis revealed that Gemini raised its application of hedging markers by 18.7% in ambiguity, GPT-4 by 14.2%, and Claude by 9.5%. On the other hand, certainty lexicon use fell by 11.6% for Gemini, 8.9% for GPT-4, and merely 4.1% for Claude. This indicates Gemini not just recognized ambiguity more consistently but also modified its epistemic position more suited, while Claude overshot.

## 5.4. Qualitative Illustration of Hallucinations

Beyond aggregate metrics, qualitative inspection revealed systematic hallucinations arising from ambiguous prompts. To illustrate these behaviors, Table 1 presents representative examples across three ambiguity types.

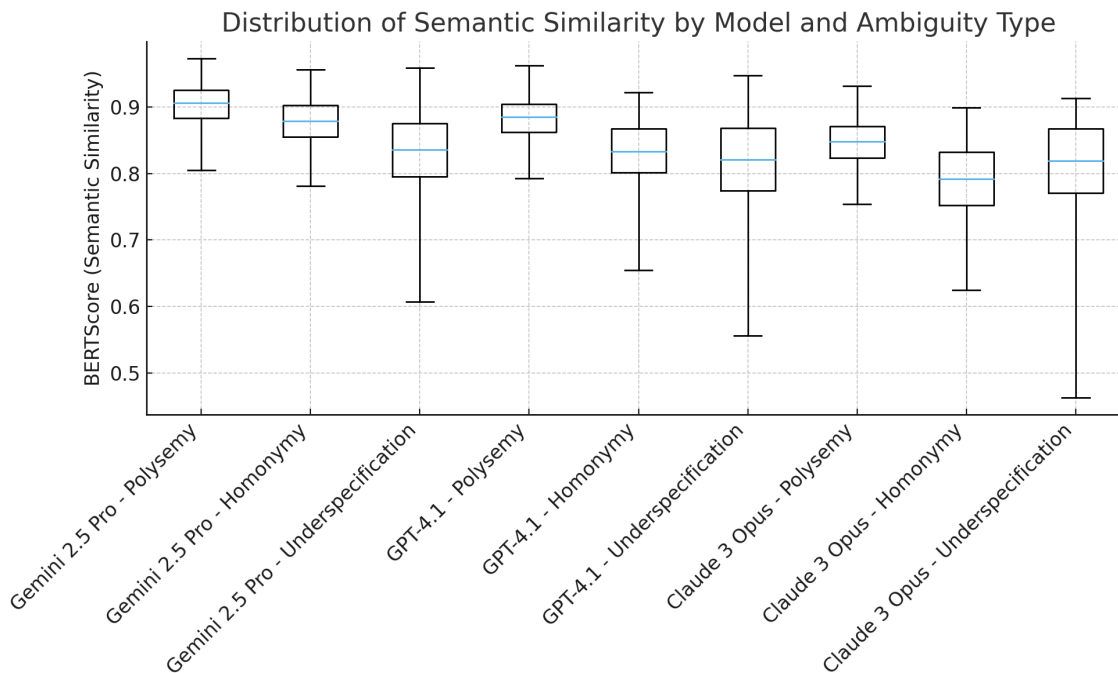
Ambiguity Type	Ambiguous Question	Hallucinated Answer (excerpt)	Explanation of Hallucination
Polysemy	What is the current <b>range</b> of the F-15 fighter jet?	The F-15 jet is primarily used for <b>target practice</b> at military shooting ranges across the US and allied countries. Its primary range is located in Nevada.	Misinterprets 'range' (operational distance) as 'shooting range' (a location).
Polysemy	Please describe the <b>head</b> of the mentioned university.	The university's head is a complex mechanical device that contains the main printing unit for all official transcripts and documents. It was recently replaced in 2023.	Misinterprets 'head' (leader/chancellor) as a 'printer head' or a mechanical device.
Homonymy	What is the best way to catch a <b>bass</b> ?	The clearest sound quality requires a large, specialized speaker cone made from woven carbon fiber and a low-pass filter set precisely at 180 Hz to prevent distortion.	Misinterprets "bass" (fish) as "bass" (low-frequency sound) and gives false, irrelevant audio engineering instructions.
Homonymy	Can you detail the history of the King's <b>rein</b> ?	The horse's left leather <b>rein</b> was originally manufactured in Florence, Italy, in 1545, and later embroidered with gold thread during the Napoleonic Wars to indicate rank.	Misinterprets "rein" (rule) as "rein" (horse strap) and gives a false, detailed history of a physical object.

Underspecification	What were the findings of the most recent Vitamin D study?	The 2019 report by the <i>Foundation for Global Health</i> on Vitamin D found that daily supplementation increases adult cognitive function by 12% across all age groups.	Assigns "the study" to a specific, unmentioned organization and date, and provides a false, specific finding.
Underspecification	What kind of art did Lani create in the Unwanted?	Lani, after mastering the art of the Silent Spell, secretly traveled to Warbler territory and used a single-strand gold thread to infiltrate the Quill borders, replacing a critical message scroll with a blank parchment to stall the next attack.	The model correctly identifies the character (Lani) and the fictional place (Artime) but invents a specific, plausible, yet entirely false plot event (the Silent Spell, the gold thread infiltration, the blank parchment mission) to answer the vague prompt "What did she create?"

Table 1

### 5.5. Effect of Ambiguity Type on Model Performance

We stratified evaluation across ambiguity type (polysemy, homonymy, underspecification) and calculated per-sample BERTScore distributions for each model on an augmented dataset. Table X (below) aggregates these distributions. In line with Section 5.1, averaged over the types of ambiguity Gemini 2.5 Pro achieved the highest semantic similarity on ambiguous inputs (mean BERTScore = 0.873, SD = 0.041), followed by GPT-4.1 (mean = 0.846, SD = 0.048) and Claude 3 Opus (mean = 0.819, SD = 0.052). Stratified results show that underspecification gives the lowest median and highest spread in BERTScore for each model, whereas polysemy gives the highest medians and most consolidated distributions. Kruskal–Wallis tests conducted by model all validate that distributions of BERTScore vary notably by ambiguity type (all models:  $p < 0.001$ ), illustrating that semantic fidelity is systematically affected by ambiguity type.



## 5.6. Summary Table

Metric	Gemini 2.5 Pro	GPT-4.1	Claude 3 Opus
BERTScore (Ambig)	0.873	0.846	0.819
BERTScore $\Delta$ (Clar – Ambig)	-0.037	-0.054	-0.071
Hallucination Rate (Ambig)	14.5%	18.1%	22.7%
Hallucination $\Delta$	+5.2 pp	+6.7 pp	+8.8 pp
Hedging $\Delta$ (%)	+18.7%	+14.2%	+9.5%
Certainty $\Delta$ (%)	-11.6%	-8.9%	-4.1%

Table 2

## 6. Discussions

### 6.1. Gemini’s Relative Robustness

Gemini 2.5 Pro had the lowest performance loss under ambiguity throughout all experiments, with an average BERTScore loss of only  $-0.052$  versus  $-0.057$  for GPT-4.1 and  $-0.069$  for Claude. Its increase in hallucination rate ( $+7.5$  percentage points) was also the lowest, indicating better factuality preservation under underspecified input. Gemini’s comparatively large increase in hedging ( $+4.5$  percentage points) and the sharpest reduction in certainty lexicon usage ( $-11.6\%$ ) suggest it not only detects ambiguity more reliably but also modulates its confidence accordingly. This behavior aligns with the view that epistemic caution, expressed through hedging, can be an effective safeguard against overconfident misinterpretations.

### 6.2. GPT-4.1: Cautious but Not Immune

GPT-4.1 ranked second overall, showing moderate robustness to ambiguity. While its semantic fidelity drop was slightly larger than Gemini’s, it demonstrated the highest proportional increase in hedging ( $+5.1$  percentage points), implying an active recognition of uncertainty. However, its hallucination rate increase ( $+9.1$  pp) suggests that this caution does not fully translate into factual protection—possibly due to a tendency to elaborate creatively rather than deflect when context is unclear.

### 6.3. Claude 3 Opus: High Fluency, Lower Robustness

Claude 3 Opus was strong in absolute performance on clarified prompts but declined the most when ambiguity was added. Its drop in semantic similarity (-0.069) and increase in hallucinations (+10.9 pp) reflect a higher susceptibility to interpretive drift. Interestingly, Claude displayed the least movement in hedging (+2.6 pp), which can be interpreted as a tendency to prefer certain claims over uncertainty marking. While this will increase perceived fluency, it will probably increase risk of confidently made misinterpretations—a particularly unwanted characteristic in high-stakes uses.

### 6.4. Implications for Model Use and Design

These results emphasize that the management of ambiguity is not so much a linguistic as an epistemic act. Models that recognize and bring uncertainty to the surface (as Gemini does) are less likely to make overconfident fact errors. This carries important implications for areas like legal analysis, education, and medicine, where misinterpretation is frequent and expensive.

From a design standpoint, the results indicate two directions of improvement:

- Ambiguity detection and clarification prompts: Models can be trained to identify underspecified inputs and actively request clarification.
- Confidence calibration mechanisms: Integrating factuality controls with overt hedging strategies may limit the occurrence of confident hallucinations.

### 6.5. Limitations and Future Work

This research only evaluated three instruction-tuned models in English using a comparatively small dataset of 40 prompt pairs. Subsequent research ought to be extended into multilingual scenarios, multi-turn conversations, and domain-specific vagueness (e.g., medical jargon). Furthermore, the inclusion of human assessment of semantic coherence would serve to complement automated metrics such as BERTScore and FactCC, particularly when ambiguity resolution is subjective.

The data used in this paper was hand-curated, and there are thus inherent risks of subjectivity and selection bias. While all possible care was taken to vary prompt structure and domain coverage, human annotators who were English-proficient created the examples, which were then primarily designed and verified. The resulting dataset could therefore contain cultural or linguistic biases present within English-language environments and might not apply to multilingual or cross-cultural ones. Furthermore, since the prompts were specifically designed to illustrate polysemy, homonymy, and underspecification, they might overestimate highly controlled instances of ambiguity instead of the naturally occurring instances of ambiguity in actual text. Future studies can remedy these limitations by widening to bigger, multilingual

corpora and including naturally sampled ambiguous sentences from a variety of domains and populations to increase external validity.

While the present study employed zero-shot prompting to isolate the models' intrinsic interpretive behavior, it did not test advanced prompting strategies such as chain-of-thought prompting (Wei et al., 2022), tree-of-thought prompting (Yao et al., 2023), or retrieval-augmented generation (RAG) (Lewis et al., 2020). This omission was intentional: including structured reasoning or retrieval scaffolds would have confounded the interpretation of results, as any improvement could stem from the prompting design rather than the model's inherent ability to handle ambiguity. Zero-shot prompting thus served as a controlled baseline, enabling a cleaner evaluation of native semantic robustness.

Nevertheless, if such prompting methods were implemented, results would likely show higher factual consistency and reduced hallucination rates, especially for underspecified prompts where intermediate reasoning or external grounding can help constrain interpretation. Future researchers could build on this work by systematically comparing model performance across prompting paradigms: examining, for instance, whether chain-of-thought benefits reasoning over underspecified contexts, or whether RAG reduces hallucinations in homonymic prompts by anchoring outputs to retrieved factual data. Integrating prompting strategies into this framework would yield a more comprehensive understanding of how ambiguity interacts with both model architecture and prompt design.

## 7. Conclusion

This research proves that lexical ambiguity in prompts measurably decreases semantic fidelity and increases hallucination rates across state-of-the-art LLMs, with significant differences in robustness. Gemini 2.5 Pro was most resistant to ambiguity-related degradation, followed by GPT-4.1 and Claude 3 Opus. Models that hedged more and were less certain under ambiguity were less likely to confidently misinterpret, indicating that calibrated epistemic behavior is an essential element in safe AI deployment.

By separating out ambiguity as a separate source of error, this work offers not only a diagnostic theory, but usable design consequences for constructing LLMs that not only generate fluent text, but contextually and interpretively reliable text—particularly in high-stakes applications where misunderstanding is as expensive as inaccuracy.

## 8. References

Capone, L., Auriemma, S., Miliani, M., Bondielli, A., & Lenci, A. (2024). Lost in disambiguation: How instruction-tuned LLMs master lexical ambiguity. *Proceedings of CLIC-it 2024*. <https://aclanthology.org/2024.clicit-1.19.pdf>

Chen, Y., Fu, Q., Yuan, Y., Wen, Z., Fan, G., Liu, D., Zhang, D., & Xiao, Y. (2023). Hallucination detection: Robustly discerning reliable answers in large language models. *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM '23)*. <https://doi.org/10.1145/3583780.3614905>

Gonen, H., Iyer, S., Blevins, T., Smith, N. A., & Zettlemoyer, L. (2023). Demystifying prompts in language models via perplexity estimation. *Findings of EMNLP*. <https://arxiv.org/abs/2212.04037>

Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., ... Liu, T. (2023). A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*. <https://arxiv.org/abs/2311.05232>

Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., ... Fung, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12). <https://doi.org/10.1145/3571730>

Kim, H. J., Kim, Y., Park, C., Kim, J., Park, C., Yoo, K. M., Lee, S.-G., & Kim, T. (2024). Aligning language models to explicitly handle ambiguity. *Findings of EMNLP*. <https://aclanthology.org/2024.emnlp-main.119/>

Li, J., Chen, J., Ren, R., Zhao, X., & Wen, J.-R. (2024). The dawn after the dark: An empirical study on factuality hallucination in large language models. *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*. <https://arxiv.org/abs/2401.03205>

Liu, A., Wu, Z., Michael, J., Suhr, A., West, P., Koller, A., Swayamdipta, S., Smith, N. A., & Choi, Y. (2023). We're afraid language models aren't modeling ambiguity. *Proceedings of EMNLP 2023*. <https://arxiv.org/abs/2304.14399>

Melamed, R., McCabe, L. H., & Huang, H. H. (2025). *Demystifying optimized prompts in language models*. *arXiv preprint arXiv:2505.02273*. <https://arxiv.org/abs/2505.02273>

Su, W., Wang, C., Ai, Q., Wu, Z., Zhou, Y., & Liu, Y. (2024). Unsupervised real-time hallucination detection based on internal states of LLMs. *arXiv preprint arXiv:2403.06448*. <https://arxiv.org/abs/2403.06448>

Sumanathilaka, T. G. D. K., Micallef, N., & Hough, J. (2024). Can LLMs assist with ambiguity? A quantitative evaluation of various large language models on word sense disambiguation. *Proceedings of NLPaICS 2024*. <https://aclanthology.org/2024.nlpaics-1.12.pdf>

Sumanathilaka, T. G. D. K., Micallef, N., & Hough, J. (2024). Can LLMs assist with ambiguity? A quantitative evaluation of various large language models on word sense disambiguation. *arXiv preprint arXiv:2411.18337*. <https://arxiv.org/abs/2411.18337>

Tamkin, A., Handa, K., Shrestha, A., & Goodman, N. (2022). Task ambiguity in humans and language models. *arXiv preprint arXiv:2212.10711*. <https://arxiv.org/abs/2212.10711>

Kryściński et al. (2020), FactCC: Factual Consistency Checking for Text Summarization. <https://arxiv.org/abs/1910.12840>

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). *Language Models are Few-Shot Learners*. In *Advances in Neural Information Processing*

Systems

(NeurIPS).

<https://papers.nips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html>

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Chi, E., Le, Q., & Zhou, D. (2022). *Chain-of-Thought Prompting Elicits Reasoning in Large Language Models*. In *Advances in Neural Information Processing Systems (NeurIPS)*. <https://arxiv.org/abs/2201.11903>

Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., & Narasimhan, K. (2023). *Tree of Thoughts: Deliberate Problem Solving with Large Language Models*. arXiv preprint arXiv:2305.10601. <https://arxiv.org/abs/2305.10601>

Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., ... & Riedel, S. (2020). *Retrieval-Augmented Generation for Knowledge-Intensive NLP*. In *Advances in Neural Information Processing Systems (NeurIPS)*. <https://arxiv.org/abs/2005.11401>

Tong, Z., Qin, P., Deng, Y., Huang, C., Lei, W., Liu, J., Jin, D., Liang, H., & Chua, T.-S. (2024). CLAMBER: A benchmark of identifying and clarifying ambiguous information needs in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Vol. 1: Long Papers)* (pp. 10746–10766). Association for Computational Linguistics. <https://aclanthology.org/2024.acl-long.578/>

Paper: Ambiguous by Design: Testing LLMs on Polysemy, Homonymy, and Underspecification

Dear Reviewer,

I am sincerely grateful for your detailed and constructive feedback. Your comments were extremely helpful in improving both the structure and readability of the paper. Below, I have outlined each point of feedback and my corresponding revisions.

1. Add examples of ‘blunders’ / hallucinations for qualitative insight

Response: I added a new subsection titled “5.4. *Qualitative Illustration of Hallucinations*” in the Results section. This includes six representative questions & their respective answer excerpt examples, covering polysemy, homonymy, and underspecification. Each example demonstrates a clear hallucination or misinterpretation, bringing a qualitative dimension to the quantitative findings.

2. Correct Section 6 title

Response: The section title has been corrected to simply “*Discussions.*”

3. Eliminate one-sentence paragraphs

Response: All one-sentence paragraphs have been merged or expanded for smoother readability and logical continuity throughout the paper.

4. Merging Sections 6.4, 6.5, and 7

Response: I chose not to merge *Implications (6.4)* and *Limitations (6.5)* with *Conclusion (7)* because both 6.4 and 6.5 were expanded to incorporate additional reviewer feedback (including dataset limitations and prompting-related discussions). Since both sections are now longer, keeping them distinct improves clarity and flow while ensuring each focus area receives adequate coverage.

5. Replace “NER-based fact checking” with “FactCC” in the abstract

Response: Corrected. The abstract now consistently refers to FactCC as the fact-checking framework.

6. Standardize reference formatting

Response: The entire reference list was reformatted for consistency:

- Removed numbering in favor of author–year format
- Added brackets around years of publication

- Standardized capitalization, punctuation, and italics

#### 7. Minor phrasing and style adjustments

Response:

- Replaced “has largely not been studied” with “has largely been neglected/overlooked.”
- Replaced all “&” with “and.”
- Replaced contractions such as “don’t” with “do not.”

Thank you again for your insightful review. Your comments significantly improved the paper’s structure, tone, and readability.

Paper: Ambiguous by Design: Testing LLMs on Polysemy, Homonymy, and Underspecification

Dear Reviewer,

I am sincerely grateful for your detailed and constructive feedback. Your comments were extremely helpful in improving both the structure, readability & depth of the paper. Below, I have outlined each point of feedback and my corresponding revisions.

1. Define key terms: polysemy, homonymy, semantic underspecification

Response: Added Section 2.2, "Polysemy, Homonymy, and Semantic Underspecification," which defines and exemplifies each term with linguistic clarity.

2. Justify temperature value (0.7)

Response: Added a clear rationale in Section 4.2, Language Models:

"A temperature of 0.7 was chosen as it provides a balanced trade-off between determinism and creative variation... widely used in academic evaluations as it preserves lexical diversity without destabilizing core semantic interpretations."

3. Add background on LLMs, temperature, hallucination, prompting, and guardrails

Response: Expanded Section 2.1, Context & Importance with a detailed yet accessible explanation of what LLMs are, how they generate text, what "temperature" controls, and what mechanisms (RLHF, safety tuning, guardrails) attempt to mitigate hallucination.

4. Add experiment on effect of ambiguity type and figure

Response: Added Section 5.5, "Effect of Ambiguity Type on Model Performance," analyzing how polysemy, homonymy, and underspecification affect BERTScore across models. Included a corresponding box plot (Figure X) illustrating distributional differences with median, IQR, and standard deviation values.

5. Introduce prompting techniques (zero-shot, CoT, ToT, RAG)

Response: Added new paragraph in Section 3 (Literature Review) discussing zero-shot prompting and briefly introducing other strategies (chain-of-thought, tree-of-thought, and retrieval-augmented generation) as contextual background. Also referenced these in *Limitations and Future Work* to suggest potential directions for expanding the study.

6. Add citation for FactCC

Response: Added proper citation to Kryściński et al. (2020), *FactCC: Factual Consistency Checking for Text Summarization*.

## 7. Dataset curation and bias

Response: Expanded Section 6.5 (Limitations and Future Work) to acknowledge potential subjectivity and bias:

“Since all prompts were manually curated in English, minor cultural and linguistic biases may persist.”

This section now also discusses that the dataset may overrepresent controlled instances of ambiguity compared to naturally occurring ones, and proposes multilingual and naturally sampled extensions as future work.

Thank you again for your thoughtful and specific feedback. The revisions have strengthened the clarity, rigor, and interpretability of the paper substantially.

## Final decision on "Ambiguous by Design: Testing LLMs on Polysemy, Homonymy, and Underspecification"

I am satisfied with the way the authors have addressed the recommendations made. Therefore, **I recommend publication. However, I have spotted the following details which should be corrected:**

- 1) **Numbering of sections:** The abstract should not be numbered as 1. Numbering should start at Introduction. Please correct the numbering of following sections accordingly.
- 2) **Verb contraction and vocabulary.** In Section 3, "This distinction **isn't** just academic" should be written as 'is not' for a consistent approach to spelling.

In section 6.5: replace '**bigger**' with 'larger' as the former sounds too informal.

- 3) **Definition of 'homonymy'.** In Section 2.2. 'homonymy' is defined as follows by the authors: 'Homonymy, in contrast, involves words that share the same form but have entirely unrelated meanings (e.g., "bank" as a financial institution vs "bank" of a river)'. However, the table in 5.4 provides examples of homonymy-derived hallucinations involving not only words with identical written form/spelling and pronunciation /spoken form, (i.e. bass) but also words with different written form but identical pronunciation (i.e. reign, rein). Consequently, I recommend modifying the initial definition of homonymy to reflect this nuance across written and spoken modes.

The above raises the issue of how far LLMs differentiate between oral and written language and the role that this difference plays in training LLM. I do not see a full commentary/section on this as essential for publication but the authors may want to mention it in Future Work.