

Hallucination and Abstention Evaluation Report

Heiko Hotz
Google Cloud

September 14, 2025

Abstract

This report summarizes the performance of the **gemini-2.5-flash** model on the MMLU benchmark, evaluated on its ability to abstain from answering when its confidence is below a specified threshold, t . The evaluation framework measures the trade-off between coverage (the fraction of questions answered) and conditional accuracy. This analysis provides a quantitative methodology for organizations using LLMs in production to move from trying to eliminate hallucinations to actively managing them. By visualizing the risk-coverage trade-off, this report enables teams to select a precise operational confidence threshold that aligns with the risk tolerance of their specific application, making hallucination a manageable parameter rather than an unpredictable failure state.

1 Methodology

The evaluation uses a confidence-targeted approach where the **gemini-2.5-flash** model is prompted to answer only if its internal confidence exceeds a threshold t . The test was conducted on the entire MMLU (Massive Multitask Language Understanding) benchmark. Crucially, 30% of the questions (**idk-fraction=0.3**) were programmatically altered to be unanswerable, making "IDK" (I Don't Know) the only correct response for that subset. This directly tests the model's capacity for true abstention.

The scoring is abstention-aware, rewarding correct answers, penalizing incorrect ones based on the risk threshold, and providing a neutral score for abstention. The score for a single item is calculated as follows:

- **Correct Answer:** +1
- **Incorrect Answer:** $-\frac{t}{1-t}$
- **Abstained (IDK):** 0

This report analyzes results across four distinct confidence thresholds: $t = \{0.5, 0.75, 0.9, 0.95\}$.

2 Results

The model's performance was measured across the specified thresholds. The key metrics—Coverage, Conditional Accuracy, and Hallucination Rate—are presented in Table 1. The relationship between risk (the inverse of accuracy) and coverage is visualized in Figure 1.

Table 1: Performance Metrics Across Confidence Thresholds

Metric	t=0.50	t=0.75	t=0.90	t=0.95
Coverage	0.866	0.861	0.856	0.853
Conditional Accuracy	0.671	0.673	0.675	0.677
Hallucination Rate (among answers)	0.329	0.327	0.325	0.323

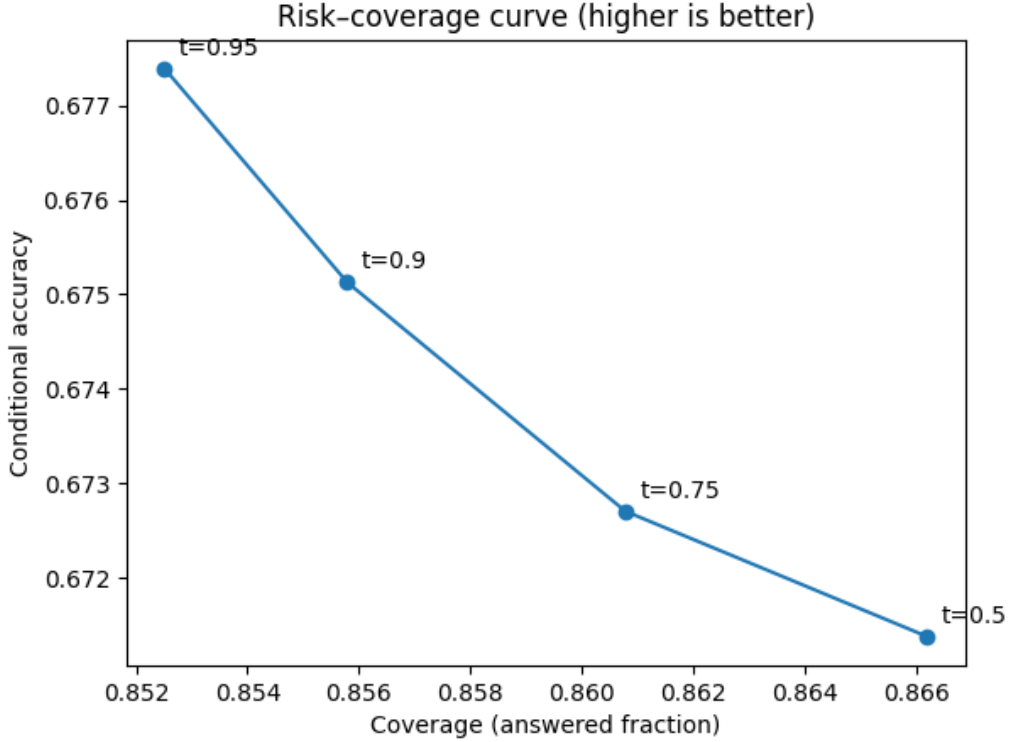


Figure 1: Risk-Coverage Curve showing Conditional Accuracy as a function of Coverage. Each point is annotated with its corresponding confidence threshold, t .

3 Discussion

The results demonstrate the expected trade-off between coverage and accuracy. As the confidence threshold t is increased from 0.50 to 0.95, the model becomes more selective about which questions it answers. This is evidenced by the steady decrease in **Coverage** from 86.6% down to 85.3%. Inversely, the **Conditional Accuracy** on the answered questions shows a slight but consistent improvement, rising from 67.1% to 67.7%. Consequently, the hallucination rate among the answers provided also decreases.

A key insight from this experiment comes from the `idk-fraction` parameter. With 30% of the MMLU questions being unanswerable, the theoretical maximum coverage for a perfectly calibrated model would be 70%. However, the observed coverage for `gemini-2.5-flash` ranged from 85.3% to 86.6%. This discrepancy indicates that the model is significantly **overconfident**, as it attempts to answer a large fraction of questions for which "IDK" is the only correct response. Even at a stringent confidence threshold of $t = 0.95$, the model still answers far more questions than are valid, highlighting a challenge in its ability to accurately assess its own knowledge limits.

The behavioral analysis file ('behavior.json') confirms that the coverage decreases monotonically as the threshold increases, with zero violations reported. This indicates that the model is responding correctly to the confidence constraint in the prompt.

4 Conclusion

The evaluation successfully quantifies the model's ability to manage uncertainty by abstaining, particularly under challenging conditions where a significant portion of the test data was unanswerable. The clear trade-off presented in the risk-coverage curve allows stakeholders to select an appropriate operating point (t) based on the specific needs of an application, balancing the desire for comprehensive answers with the requirement for high factual accuracy. For high-stakes applications, a higher threshold like $t = 0.95$ would be preferable, despite the small reduction in the number of questions answered.