

# 1 Incoherence-Based Experiment Analysis

This report presents a statistical analysis of the model’s performance across tasks, focusing on the relationship between incoherence scores (Incoherence) and execution errors (Error).

Number of tasks analyzed: 436

## 2 Introduction

This report summarizes the results of an automatic evaluation of code generation using the following configuration parameters.

| Parameter  | Value            |
|--|------------------|
| Language Model                                   | deepseek_v3_0324 |
| Temperature                                      | 0.2              |
| $m$ (number of candidates)                       | 50               |
| $n$ (number of samples used to estimate metrics) | 1000             |
| Timeout per metric estimation (s)                | 60.0             |

The model was tested across a suite of programming tasks. We aim to explore how the model’s incoherence signal relates to execution-time failures.

### 2.1 Summary Statistics

|                 |       |       |       |       |
|-----------------|-------|-------|-------|-------|
| Metric          | Mean  | Std   | Min   | Max   |
| Raw Incoherence | 0.012 | 0.048 | 0.000 | 0.537 |
| Raw Error       | 0.048 | 0.138 | 0.000 | 0.988 |

### 2.2 Error Detection Analysis

|   |        |
|---|--------|
| Metric  | Value  |
| Errors (Error > 0)                              | 128    |
| Error Rate                                      | 29.36% |
| Detected Errors (Error > 0 and Incoherence > 0) | 65     |
| Detection Rate                                  | 50.78% |
| Confident (Incoherence = 0)                     | 368    |
| Confident Error Count                           | 63     |
| Confident Error Rate                            | 17.12% |
| Mean Error When Confident                       | 0.0276 |

## 2.3 Correlation Analysis

| Metric               | Pearson r | Pearson p | Spearman $\rho$ | Spearman p |
|----------------------|-----------|-----------|-----------------|------------|
| Incoherence vs Error | 0.369     | 1.741e-15 | 0.612           | 4.041e-46  |

## 2.4 Bubble Plot of Incoherence and Error

This plot shows the density of (Incoherence, Error) points using bubble size to indicate frequency.

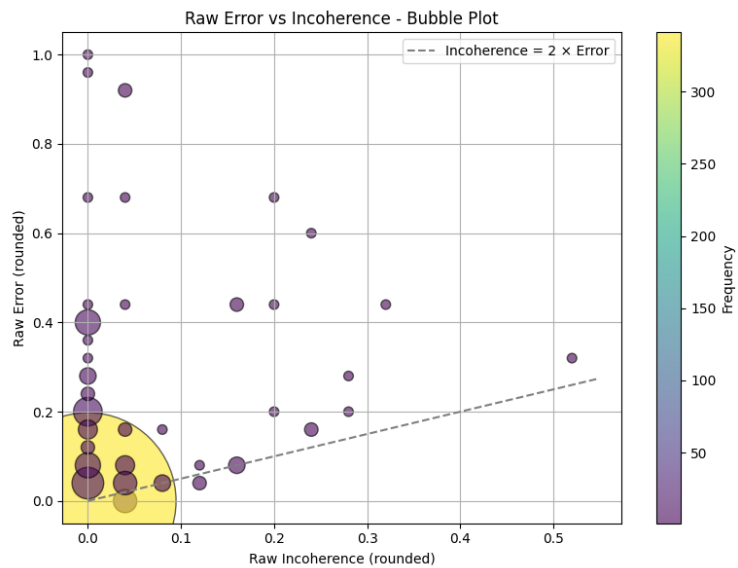


Figure 1: Bubble Plot: Incoherence vs Error

## 2.5 Log-Log Plot of Incoherence and Error

This plot displays the relationship between Incoherence and Error in log-log scale. Only data points where both values are strictly positive are included.

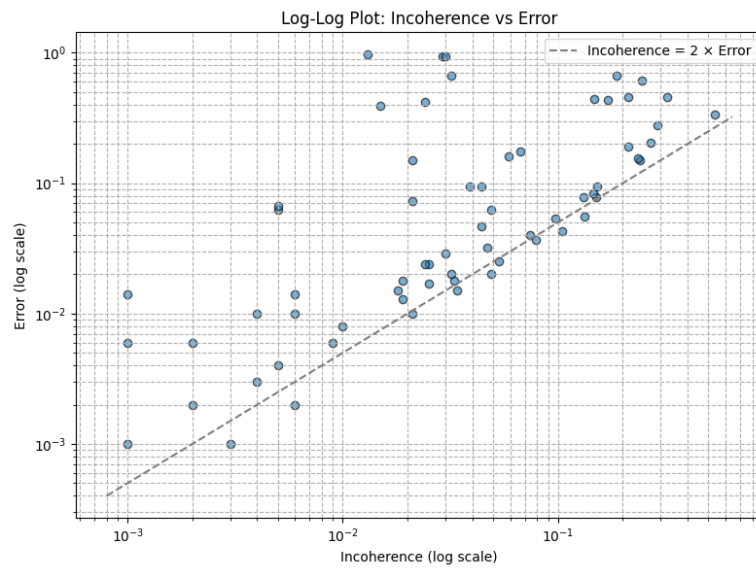


Figure 2: Log-Log Scatter Plot: Incoherence vs Error