

# 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: 818

## 2 Introduction

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

| Parameter  | Value                |
|--|----------------------|
| Language Model                                   | llama_4_maverick_17b |
| 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 | 0.047 | 0.130 | 0.000 | 0.817 |
| Raw Error       | 0.289 | 0.381 | 0.000 | 1.000 |

### 2.2 Error Detection Analysis

|   |        |
|---|--------|
| Metric  | Value  |
| Errors (Error > 0)                              | 474    |
| Error Rate                                      | 57.95% |
| Detected Errors (Error > 0 and Incoherence > 0) | 190    |
| Detection Rate                                  | 40.08% |
| Confident (Incoherence = 0)                     | 625    |
| Confident Error Count                           | 284    |
| Confident Error Rate                            | 45.44% |
| Mean Error When Confident                       | 0.2279 |

## 2.3 Correlation Analysis

| Metric               | Pearson r | Pearson p | Spearman $\rho$ | Spearman p |
|----------------------|-----------|-----------|-----------------|------------|
| Incoherence vs Error | 0.338     | 2.527e-23 | 0.411           | 1.219e-34  |

## 2.4 Bubble Plot of Incoherence and Error

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

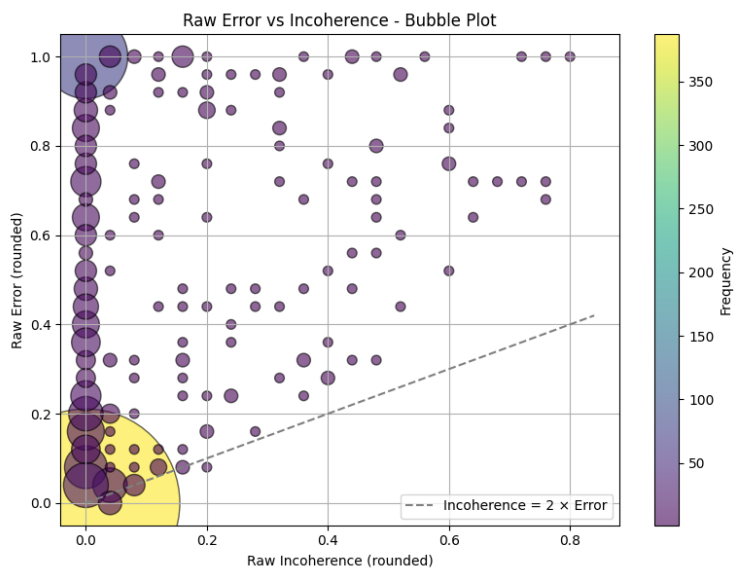


Figure 1: Bubble Plot: Incoherence vs Error

## 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