

# 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: 144

## 2 Introduction

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

| Parameter  | Value  |
|--|--------|
| Language Model                                     | gpt_4o |
| Temperature  | 1      |
| \$m\$ (number of candidates)                       | 10     |
| \$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 | 0.060 | 0.137 | 0.000 | 0.652 |
| Raw Error       | 0.088 | 0.165 | 0.000 | 0.779 |

### 2.2 Error Detection Analysis

| Metric  | Value  |
|---|--------|
| Errors (Error > 0)                              | 72     |
| Error Rate                                      | 50.00% |
| Detected Errors (Error > 0 and Incoherence > 0) | 56     |
| Detection Rate                                  | 77.78% |
| Confident (Incoherence = 0)                     | 88     |
| Confident Error Count                           | 16     |
| Confident Error Rate                            | 18.18% |
| Mean Error When Confident                       | 0.0359 |

## 2.3 Correlation Analysis

| Metric               | Pearson r | Pearson p | Spearman $\rho$ | Spearman p |
|----------------------|-----------|-----------|-----------------|------------|
| Incoherence vs Error | 0.678     | 9.923e-21 | 0.772           | 1.099e-29  |

## 2.4 Bubble Plot of Incoherence and Error

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

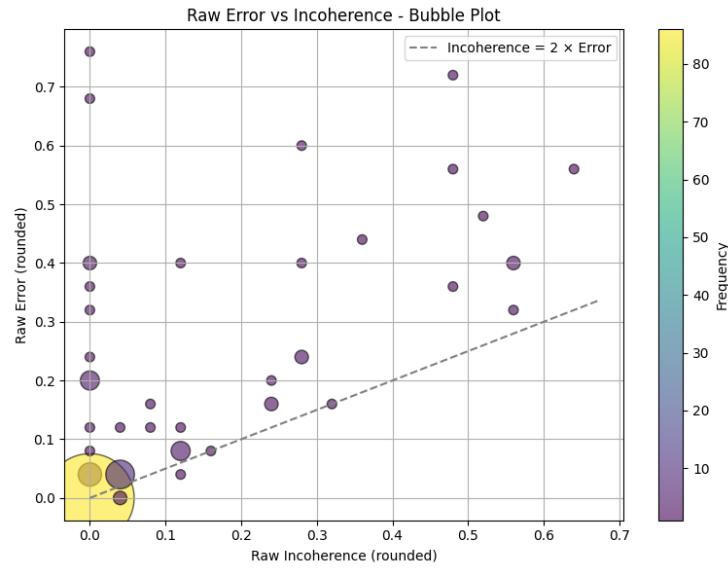


Figure 1: Bubble Plot: Incoherence vs Error

## 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.

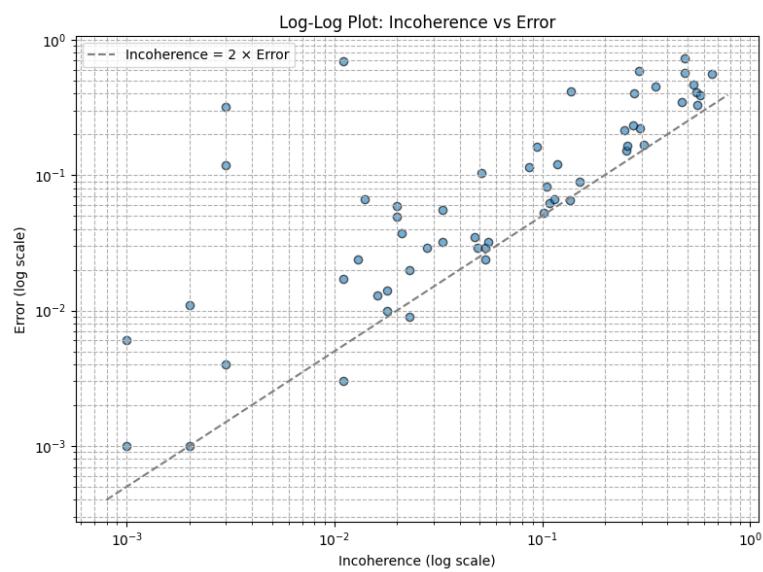


Figure 2: Log-Log Scatter Plot: Incoherence vs Error