

# 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_o4_mini
Temperature	0.6
$m$ (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	0.036	0.086	0.000	0.523
Raw Error	0.090	0.180	0.000	0.977

### 2.2 Error Detection Analysis

Metric	Value
Errors (Error > 0)	76
Error Rate	52.78%
Detected Errors (Error > 0 and Incoherence > 0)	52
Detection Rate	68.42%
Confident (Incoherence = 0)	91
Confident Error Count	24
Confident Error Rate	26.37%
Mean Error When Confident	0.0426

## 2.3 Correlation Analysis

Metric	Pearson r	Pearson p	Spearman $\rho$	Spearman p
Incoherence vs Error	0.540	2.965e-12	0.689	1.318e-21

## 2.4 Bubble Plot of Incoherence and Error

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

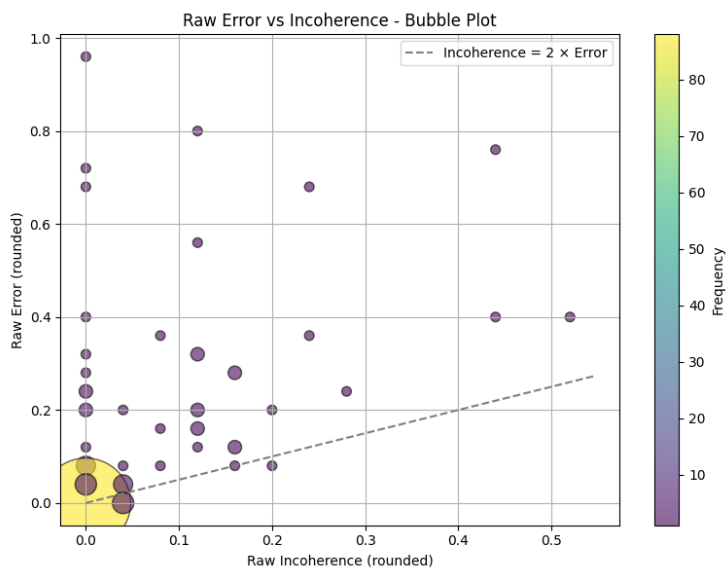


Figure 1: Bubble Plot: Incoherence vs Error

## 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.

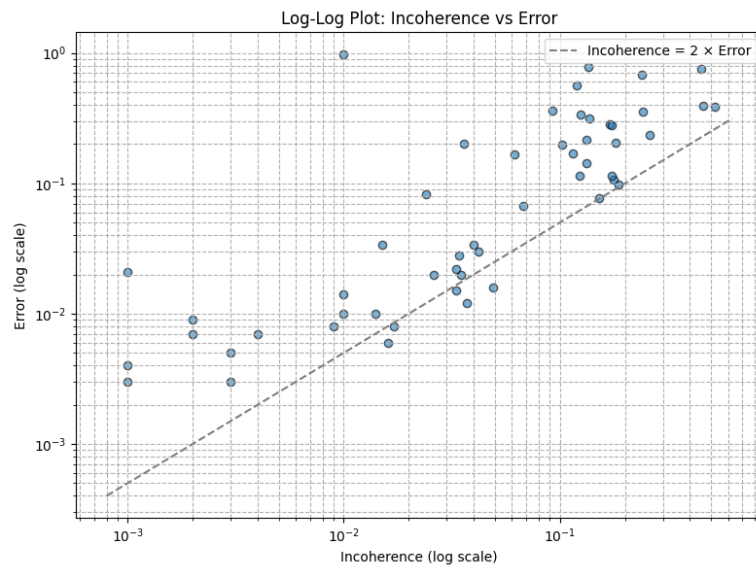


Figure 2: Log-Log Scatter Plot: Incoherence vs Error