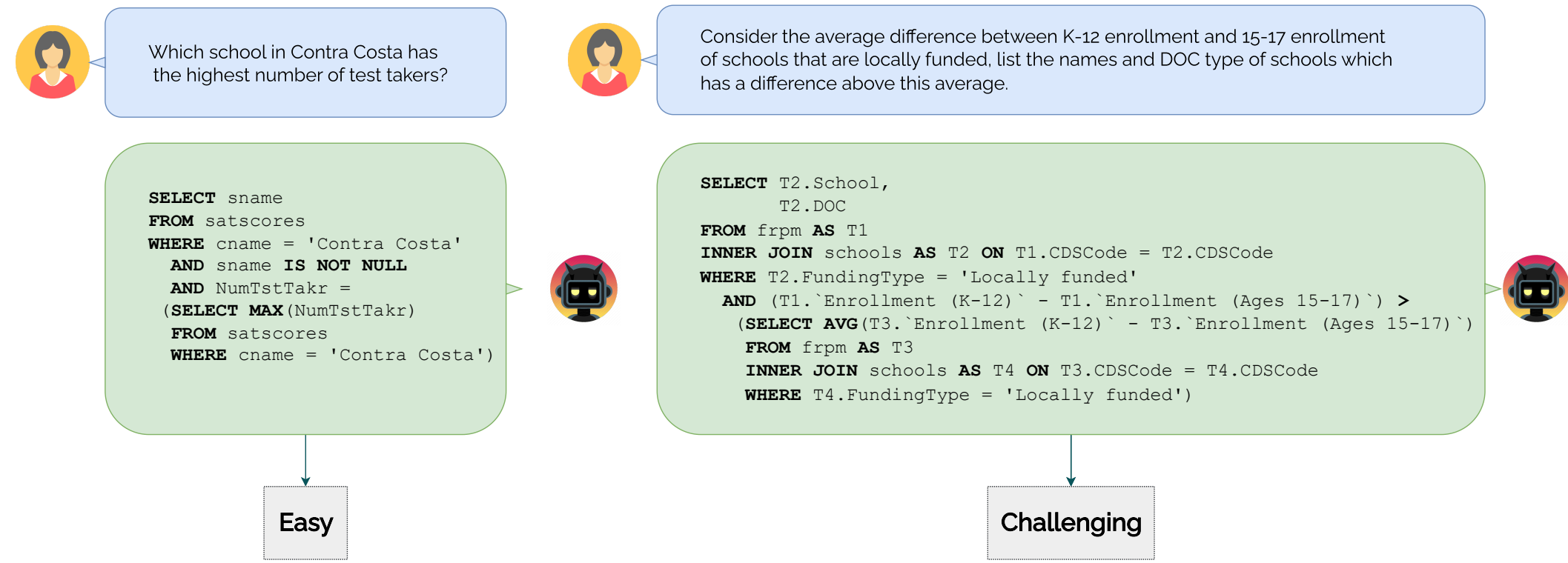# Given an NL query, route to the cheapest LLM capable of generating executable SQL
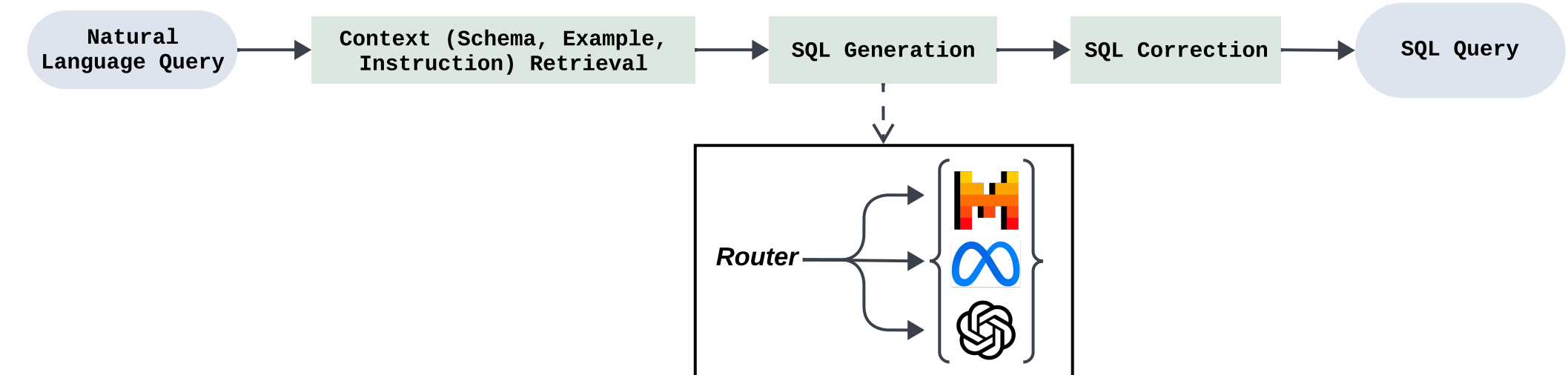
## Introduction

**Text-to-SQL** allows users to query databases in natural language, making data more accessible to non-experts. However, using powerful LLMs like `GPT-4o` for all queries leads to unnecessary costs and latency, especially for simpler queries.



## Problem statement

Given a set of $N$ models $\{\mathcal{M}_0, \mathcal{M}_1, \ldots, \mathcal{M}_{N-1}\}$ with varying SQL generation capabilities and costs, the goal is to select the weakest model $\mathcal{M}_i$ that can generate accurate SQL for a query $Q$, balancing accuracy, cost, and latency. We use a dataset $H$ of past queries and SQL outputs to create an $N$-ary routing function that assigns a query to the best model $\mathcal{M}_l$ (where $l \in [0, N-1]$), or determines that no model can generate SQL ($l = N$).



## Datasets & Metrics



**Dataset:** BIRD →1534 dev queries and 9428 training queries

**Metric:** Execution Accuracy ($EX$) measures the proportion of queries in the evaluation set $S$ where the predicted SQL's output matches the ground-truth relation, with $0 \leq EX(S) \leq 1$.

$$EX = \frac{1}{N} \sum_{n=1}^{N} \mathbf{1}(V_n, \hat{V}_n) \quad \text{where} \quad \mathbf{1}(V_n, \hat{V}_n) = \begin{cases} 1 & \text{if } V_n = \hat{V}_n \\ 0 & \text{if } V_n \neq \hat{V}_n \end{cases}$$

## Failure Analysis

The analysis shows a clear gap in model capabilities, with high overlap in failed queries. This suggests that routing primarily reduces costs, without improving $EX$ beyond the strongest model.
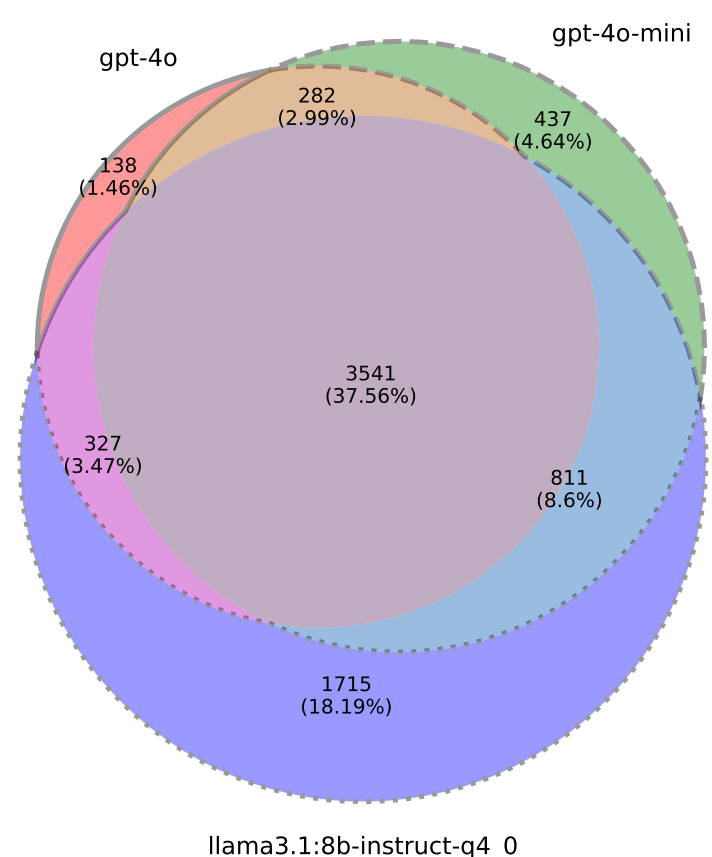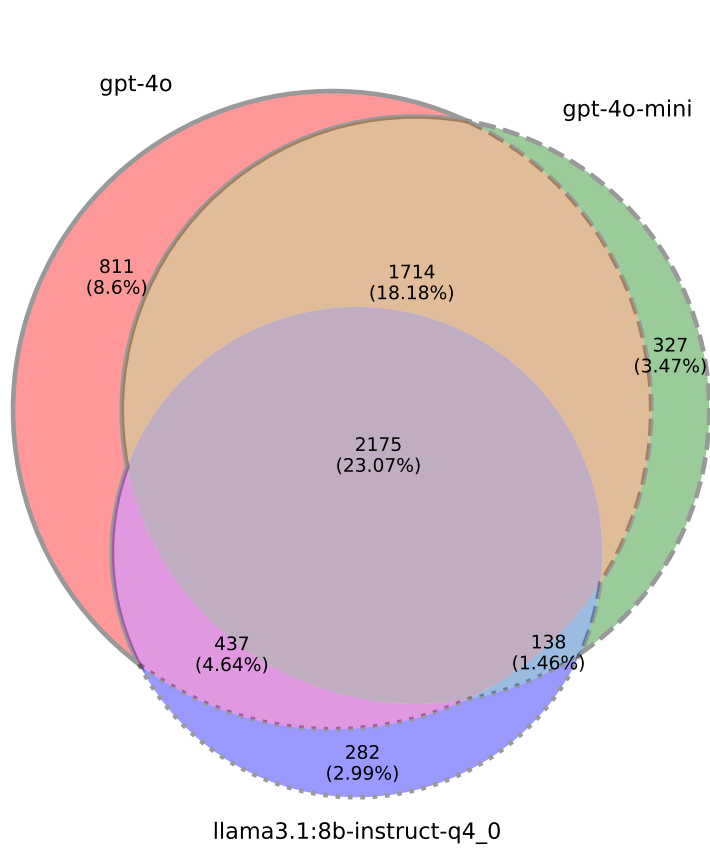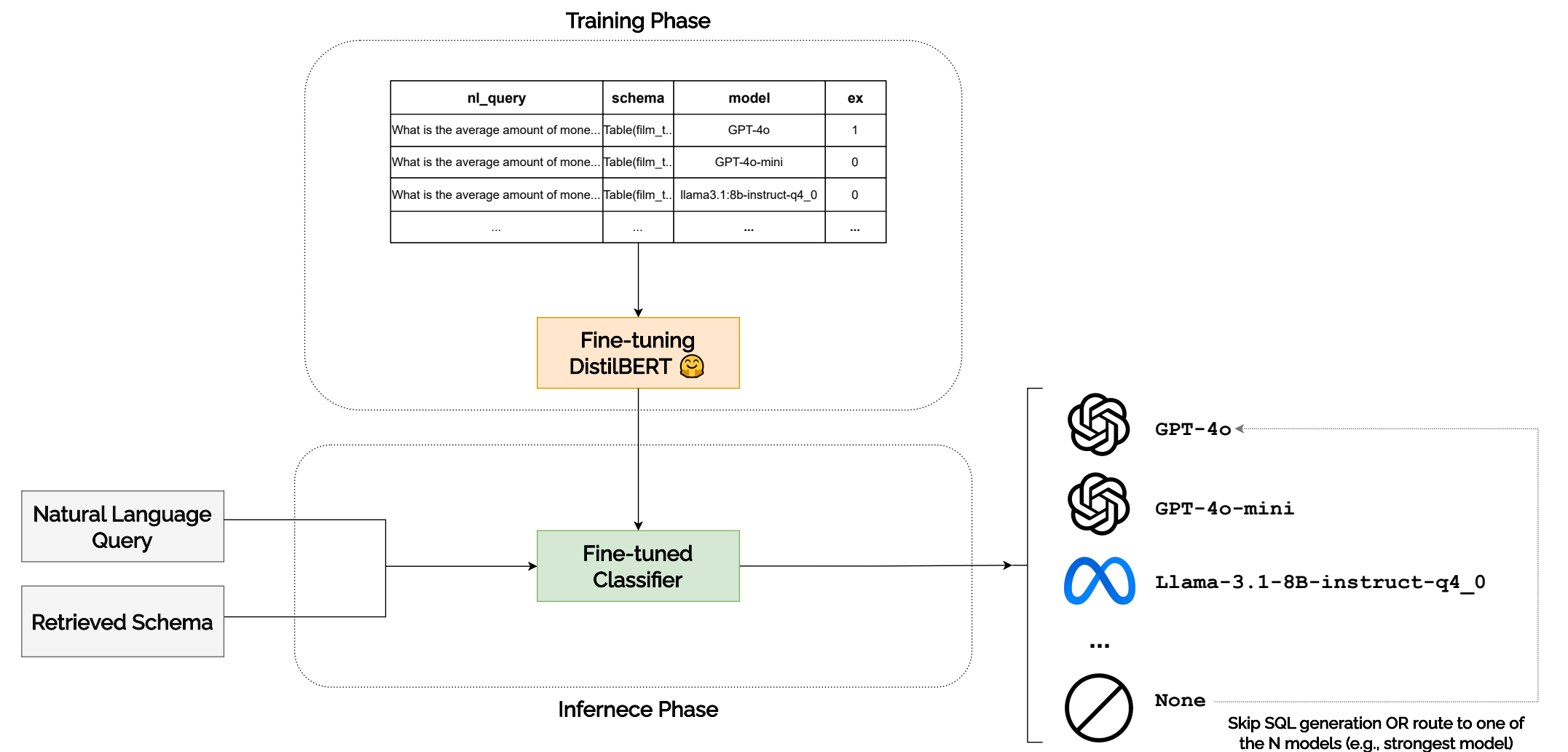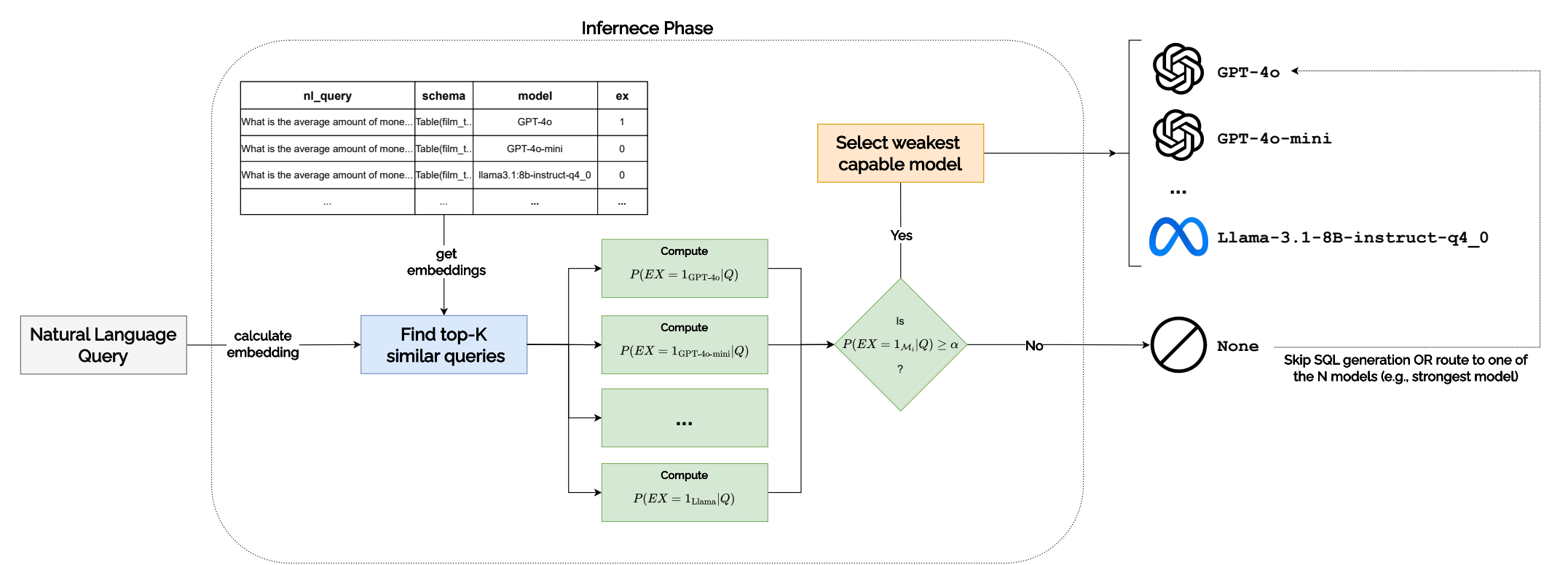


Figure 1. Failure Cases Dist.



Figure 2. Correct Cases Dist.

## Methodology

### Classification-Based Routing ($R_{BERT}$)



### Score-Based Routing ($R_k^\alpha$)



## Results

| Gen. | EX% | gpt-4o | 4o-mini | Llama | None | Cost Red. |
|------|-----|--------|---------|-------|------|-----------|
| `gpt-4o` | 61.02 | 1534 | – | – | – | 1x |
| `4o-mini` | 49.22 | – | 1534 | – | – | 16.6x |
| `Llama` | 29.34 | – | – | 1534 | – | ∞ |
| $R_{25}^{0.7}$ | 60.14 | 197 | 88 | 5 | 1243 | 1.1x |
| $R_{10}^{0.8}$ | 59.42 | 160 | 127 | 26 | 1220 | 1.1x |
| $R_{24}^{0.6}$ | 57.92 | 324 | 265 | 54 | 890 | 1.3x |
| $R_{BERT}$ | 55.21 | 118 | 311 | 167 | 938 | 1.4x |

## Accuracy & Cost Trade-Off

In the score-based approach, by adjusting $K$ (similar queries) and $\alpha$ (threshold), we balance $EX$ and cost. Higher values improve $EX$ but rely more on the strongest model, increasing cost.
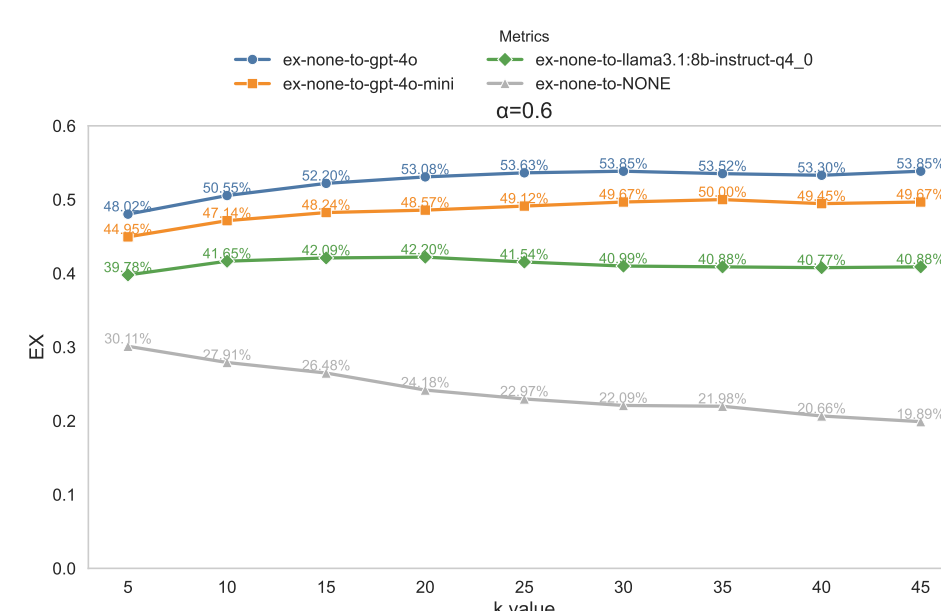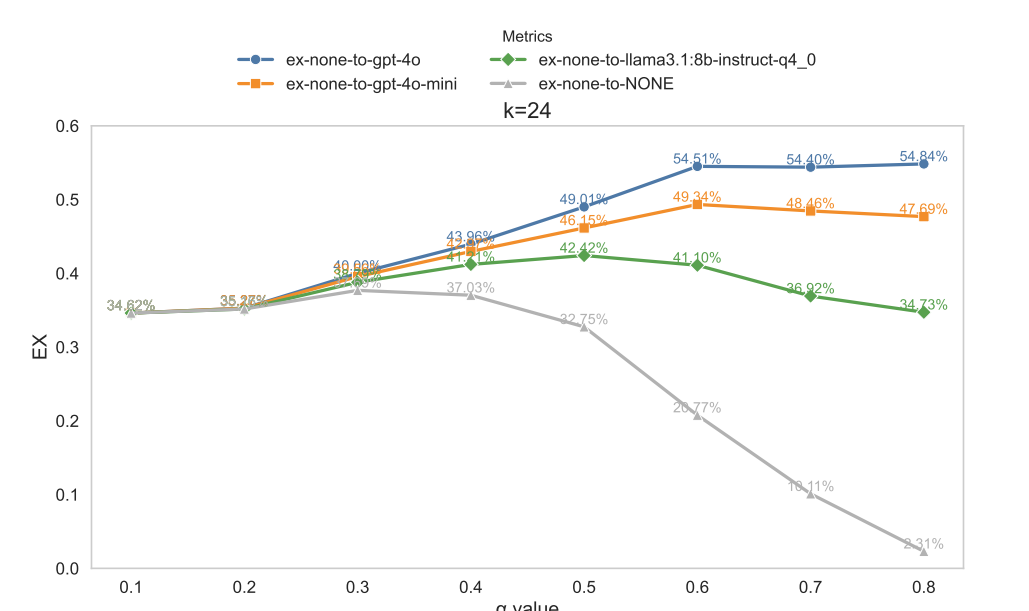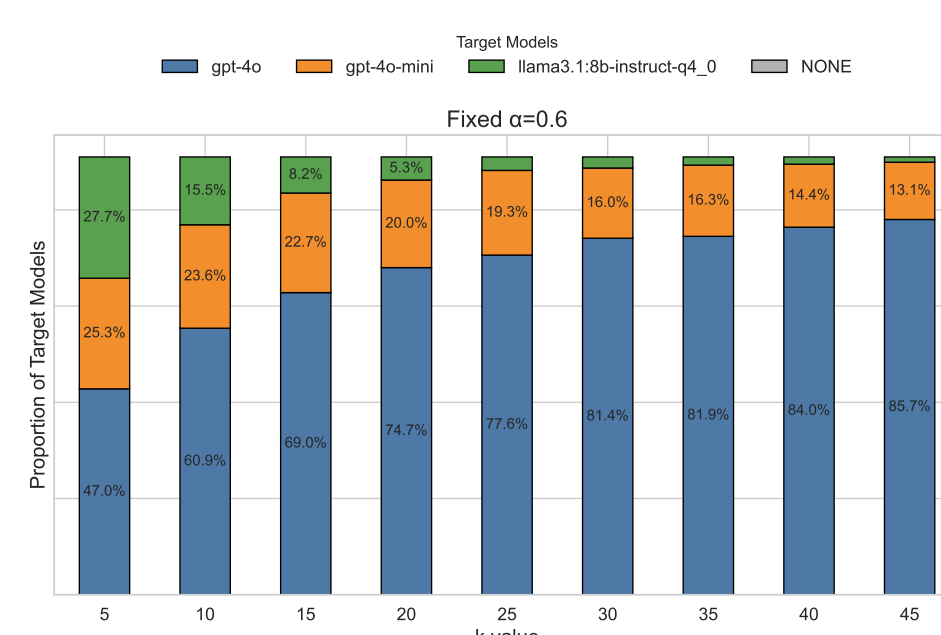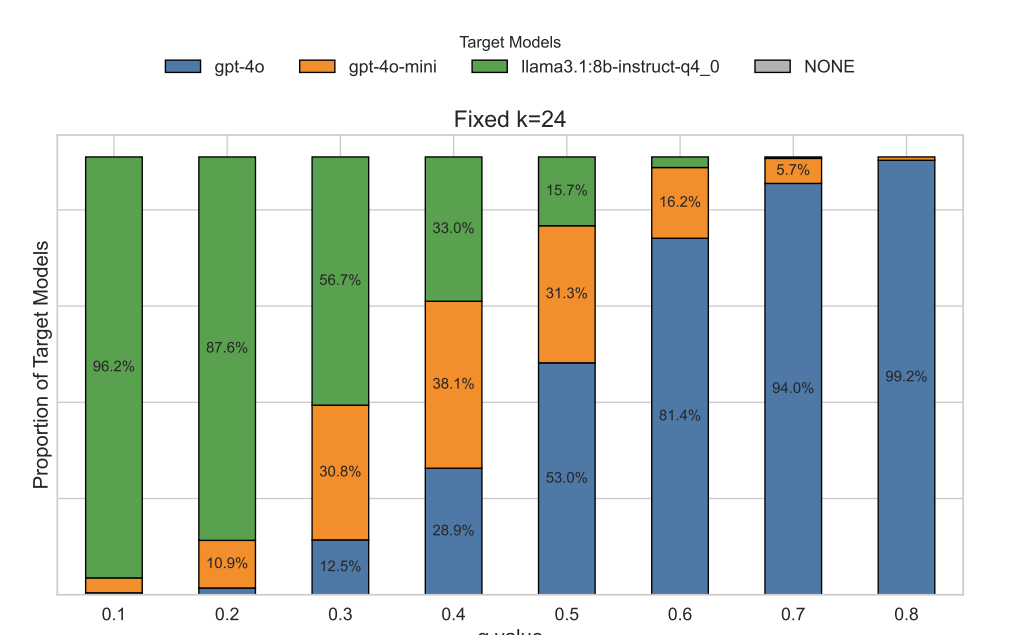


Figure 3. $EX$ vs. $K$



Figure 4. $EX$ vs. $\alpha$



Figure 5. Dist. vs. $K$ (None→`gpt-4o`)



Figure 6. Dist. vs. $\alpha$ (None→`gpt-4o`)