

# INFERENCE DYNAMICS: Adaptive LLM Routing through Structured Capability and Knowledge Profiling

Anonymous ACL submission

## Abstract

Large Language Model (LLM) routing is a pivotal technique for navigating a diverse landscape of LLMs, enabling the selection of the best-performing LLMs for specific user queries while balancing performance and cost. However, current routing approaches often face limitations in scalability when dealing with a large pool of specialized LLMs, or in their adaptability to extending model scope and evolving capability domains. To overcome those challenges, we propose **InferenceDynamics**, a flexible and scalable multi-dimensional routing framework by modeling the capability and knowledge of models. We operate it on our comprehensive dataset **RouteMix**, and demonstrate its effectiveness and generalizability in group-level routing using modern benchmarks including MMLU-Pro, GPQA, BigGenBench, and LiveBench, showcasing its ability to identify and leverage top-performing models for given tasks, leading to superior outcomes with cost efficiency. The broader adoption of InferenceDynamics can empower users to harness the full specialized potential of the LLM ecosystem, and our code will be made publicly available to encourage further research.

## 1 Introduction

The rapid proliferation of Large Language Models (LLMs) has unveiled a rich landscape of specialized capabilities, with different models demonstrating unique strengths across a multitude of domains and tasks (Matarazzo and Torlone, 2025; Li et al., 2024a). This specialization necessitates a sophisticated approach to model selection, where the primary goal is to identify and utilize the LLM best suited to the specific demands of a user’s query. LLM routing (Chen et al., 2025) emerges as a critical paradigm to address this, creating mechanisms to strategically dispatch queries to optimal model from a diverse pool, thereby maximizing performance, relevance, and the quality of outcomes,

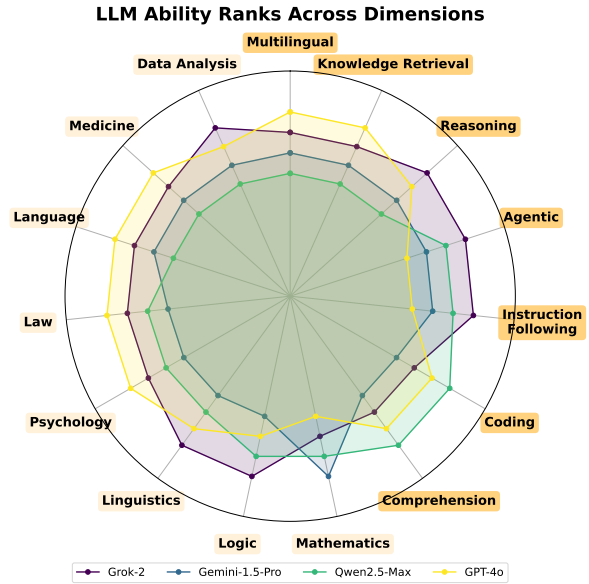


Figure 1: Quantification of Knowledge and Capability of top 4 models among candidate LLMs.

while also considering factors like inference cost and latency.

Early explorations in LLM routing often simplified the selection problem, for instance, by framing it as a binary classification task—e.g., choosing between a generalist small model and a powerful large model. Methods such as AutoMix (Aggarwal et al., 2024), HybridLLM (Ding et al., 2024), and RouteLLM (Ong et al., 2025) have demonstrated the effectiveness of this approach, primarily emphasizing the trade-off between cost and performance. Moreover, the emergence of advanced models like GPT-5 (OpenAI, 2025) further substantiates its validity. While valuable for two-model scenarios, such binary frameworks face inherent scalability challenges, as selecting the optimal model from many candidates using only pairwise comparisons becomes computationally costly and inefficient.

More recent works have advanced the field by leveraging richer model representations to better

062 evaluate and route LLMs based on their specific  
063 capabilities. While methods including RouterDC  
064 (Chen et al., 2024), C2MAB-V (Dai et al., 2024),  
065 and P2L (Frick et al., 2025) offer more sophisti-  
066 cated mechanisms for capturing model strengths,  
067 their primary limitation lies in the significant re-  
068 training or recalibration required to effectively  
069 support newly introduced LLMs, hindering their  
070 agility in a rapidly evolving model landscape.  
071 Model-SAT (Zhang et al., 2025) addresses this lim-  
072 itation with human-specified, model-agnostic de-  
073 compositions of query knowledge. However, its  
074 reliance on static knowledge sets limits adaptabil-  
075 ity to new knowledge dimensions and hinders fine-  
076 grained evaluation in specialized domains. More-  
077 over, embedding model-specific knowledge into  
078 prompts causes prompt lengths to grow rapidly  
079 with the number of models and knowledge.

080 To address this gap, we introduce **Inference-**  
081 **Dynamics**, a novel system designed for perfor-  
082 mant, scalable, and adaptable LLM routing. **In-**  
083 **ferenceDynamics** operates by extracting capabil-  
084 ity requirements and domain-specific knowledge  
085 from incoming queries, modeling the correspond-  
086 ing capabilities and knowledge profiles of avail-  
087 able LLMs, and then intelligently routing queries  
088 to the most suitable models. To demonstrate  
089 the effectiveness and generalizability of our ap-  
090 proach, we constructed a comprehensive dataset  
091 aggregated from 24 diverse benchmarks. We  
092 then evaluated our routing algorithm on four chal-  
093 lenging out-of-distribution (OOD) benchmarks:  
094 MMLU-Pro (Wang et al., 2024b), GPQA (Rein  
095 et al., 2023), BigGenBench (Kim et al., 2024),  
096 and LiveBench (White et al., 2025). Experimental  
097 results show that our routing algorithm achieved  
098 the highest average score, surpassing the top-  
099 performing single LLM by a substantial margin of  
100 1.22 points under optimal routing conditions. Fur-  
101 thermore, when operating under cost constraints,  
102 our algorithm delivered competitive performance  
103 comparable to the best single LLM, while utilizing  
104 nearly half the budget.

105 The contributions of our work are summarized  
106 as follows:

- 107 • We introduce **RouteMix**, a comprehensive  
108 dataset aggregated from 24 diverse bench-  
109 marks, specifically curated for rigorously eval-  
110 uating the generalization capabilities of LLM  
111 routing algorithms.
- 112 • We propose **InferenceDynamics**, an efficient

113 routing algorithm demonstrating generaliza-  
114 tion capabilities on previously unseen queries.

- 115 • Experimental results validate that **Inference-**  
116 **Dynamics** significantly enhances LLM rout-  
117 ing, substantially outperforming the leading  
118 single model while concurrently reducing  
119 cost.

## 120 2 Related Works

### 121 2.1 Multi-LLM System

122 A Multi-LLM system (Chen et al., 2025) refers to  
123 the architecture that combines LLMs to collabora-  
124 tively solve tasks more effectively than any single  
125 model. The rapid proliferation of diverse LLMs  
126 has spurred significant interest in such systems,  
127 which are realized through several architectural  
128 patterns. LLM ensembling (Jiang et al., 2023; Li  
129 et al., 2024b) enhances accuracy or robustness by  
130 processing the same input through several models  
131 and then aggregating their responses. Cascaded  
132 systems (Zhang et al., 2024; Kolawole et al., 2024;  
133 Chen et al., 2023) strategically employ a sequence  
134 of models—often initiating with smaller, faster  
135 LLMs for initial processing or simpler queries and  
136 escalating to more powerful, resource-intensive  
137 ones only when necessary—thereby optimizing  
138 resource use. Furthermore, the development of  
139 collaborative LLM agents (Wang et al., 2024a; Xu  
140 et al., 2024, 2025) involves multiple LLMs, with  
141 distinct roles or access to different tools, interact-  
142 ing to address complex, multi-step problems that  
143 demand sophisticated coordination. While these  
144 multi-LLM approaches demonstrate considerable  
145 advancements, they often necessitate querying mul-  
146 tiple models, which can increase computational  
147 cost and latency. Moreover, as the number and  
148 diversity of available LLMs continue to grow, it  
149 becomes critical to route queries to the most suit-  
150 able model, effectively balancing performance with  
151 operational costs.

### 152 2.2 LLM Routing

153 LLM routing seeks to identify the most suitable  
154 language model for a given query, with various  
155 strategies proposed. Early methods include LLM-  
156 Blender (Jiang et al., 2023), which employs an en-  
157 semble framework querying multiple LLMs to se-  
158 lect the optimal response, and AutoMix (Aggarwal  
159 et al., 2024), which utilizes a smaller model for self-  
160 verification before potentially escalating to a larger  
161 model. While these can improve performance,

their reliance on multiple querying inherently increases latency. Other strategies, such as HybridLLM (Ding et al., 2024) and RouteLLM (Ong et al., 2025), focus on training a binary classifier to choose between a human-defined strong and weak model. However, these methods’ efficacy is highly contingent on the subjective definition of model strength and can be computationally expensive when applied to a large pool of LLMs. More recent research has shifted towards multi-LLM routing. RouterDC (Chen et al., 2024), C2MAB-V (Dai et al., 2024), and Prompt-to-Leaderboard (Frick et al., 2025) trains a parametric router to route queries. Concurrently, ModelSpider (Zhang et al., 2023) and EmbedLLM (Zhuang et al., 2025) encode LLMs into learnable representations to facilitate routing. Despite these advancements, a significant limitation is the need to retrain the entire routing mechanism when new models are introduced. Addressing this, Model-SAT (Zhang et al., 2025) aimed to resolve the retraining weakness through human-defined, model-independent capability decompositions. However, its reliance on predefined knowledge sets limits adaptability to new dimensions, and the approach substantially increases input token counts as the number of models and knowledge components grows.

### 3 Methodology

In this section, we introduce **InferenceDynamics**, which involves: (i) identifying the knowledge and capability required for a given query, (ii) quantifying the knowledge and capability of LLMs, and (iii) routing queries to LLMs based on their scores.

#### 3.1 Problem Setup

Let  $\mathcal{M}_T = \{M_1, M_2, \dots, M_t\}$  denote a set of LLMs, and let  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_n$  be a dataset where  $\mathbf{x}_i$  represents a query and  $y_i$  its corresponding ground truth. For an unseen query  $x \in \mathcal{Q}$ , where  $x \notin \mathcal{D}$ , LLM routing is formalized as a function  $\mathcal{R} : \mathcal{Q} \rightarrow \mathcal{M}_T$ . This function maps the query  $x$  to the model  $M_{\text{best}} \in \mathcal{M}_T$  that is considered most suitable, based on a joint assessment of both cost and performance. Our objective is to develop a routing algorithm with the dataset  $\mathcal{D}$ , that effectively generalizes to OOD queries.

#### 3.2 Knowledge and Capability Generation

It is widely acknowledged that no single LLM demonstrates universal proficiency across the full spectrum of query types. Previous research (Wang

et al., 2024c; Li et al., 2024c) substantiates that distinct queries necessitate specific underlying capabilities and domain-specific knowledge. Accordingly, assessing an LLM’s aptitude for a given query necessitates identifying the requisite capabilities and knowledge pertinent to that query. Let  $\mathcal{C}$  denote the set of defined LLM capabilities and  $\mathcal{K}$  represent the world knowledge space. For a given query  $x$ , we utilize an auxiliary LLM  $\mathcal{M} \notin \mathcal{M}_T$  to predict two sets:  $\mathcal{C}_x = \{c_1, c_2, \dots \mid c_i \in \mathcal{C}\}$ : This set comprises the capabilities deemed necessary to address query  $x$ , ranked in descending order of importance.  $\mathcal{K}_x = \{k_1, k_2, \dots \mid k_i \in \mathcal{K}\}$ : This set encompasses the knowledge areas considered essential for resolving query  $x$ , also ranked in descending order of importance.

#### 3.3 Scoring

To quantify the proficiency of a model  $M_t$  with respect to specific capabilities and knowledge, we utilize the accessible set  $\mathcal{D}$ . The performance score  $s_i^t$  of model  $M_t$  for a given query-response pair  $(\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{index}}$  is determined by averaging over  $K$  independent trials:

$$s_i^t = \frac{1}{K} \sum_{k=1}^K \text{eval}(M_t(\mathbf{x}_i)_k, y_i)$$

where  $M_t(\mathbf{x}_i)_k$  is the model’s  $k$ -th generated response to the input query  $\mathbf{x}_i$ , and  $\text{eval}(\cdot, \cdot)$  represents the query-specific evaluation metric employed to compare the model’s response against the ground truth  $y_i$ . To incorporate the trade-off between performance and computational expenditure, we record the average computational cost  $\mathbf{c}_i^t$  incurred by model  $M_t$  when processing query  $\mathbf{x}_i$ .

Subsequent to the identification of the knowledge and capability sets and computing the scores for all queries in the set  $\mathcal{D}$ , we define a refined score for model  $M_t$ . This score,  $S_\beta^\alpha(M_t, \mathbf{x}_i, e)$ , quantifies the model’s effectiveness for a specific element  $e$  (which can be a knowledge item  $k \in \mathcal{K}_{\mathbf{x}_i}$  or a capability  $c \in \mathcal{C}_{\mathbf{x}_i}$ ) associated with query  $\mathbf{x}_i$ . Illustrating with a knowledge element  $k$ , this score is formulated as:

$$S_\beta^\alpha(M_t, \mathbf{x}_i, k) = \sum_{j=1}^{|\mathcal{K}_{\mathbf{x}_i}|} (s_i^t - \beta \mathbf{c}_i^t) \mathbb{1}(k = k_j) \frac{\alpha^{j-1}}{\sum_{m=1}^{|\mathcal{K}_{\mathbf{x}_i}|} \alpha^{m-1}}$$

In this formulation, the hyperparameter  $\alpha$  serves to attenuate the influence of less critical knowledge

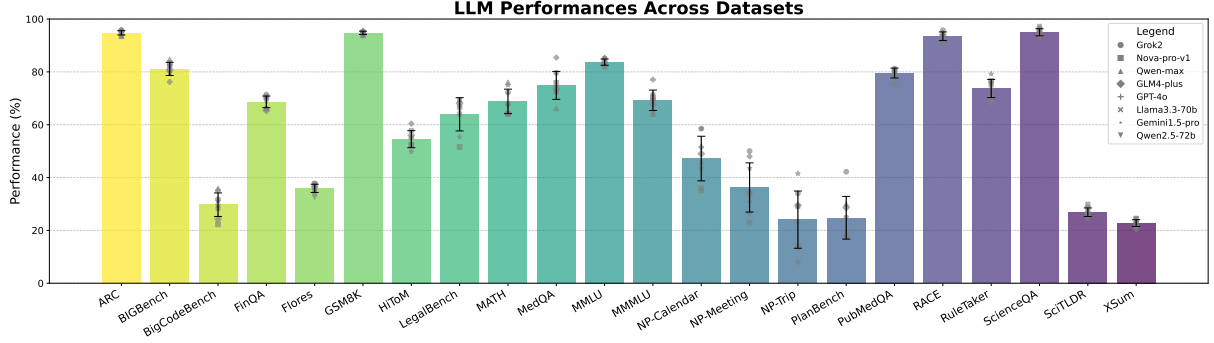


Figure 2: LLM performances across 20 datasets in **RouteMix**. Dataset labels including "PlanBench" indicate subsets of the PlanBench benchmark. For detailed metric information, refer to [Appx. §A](#).

elements, based on their rank  $j$ . The hyperparameter  $\beta$  acts as a coefficient penalizing higher computational costs. The denominator,  $\sum_{k=1}^{|\mathcal{K}_{\mathbf{x}_i}|} \alpha^{k-1}$ , functions as a normalization factor, ensuring that each query contributes equitably to the knowledge score, regardless of the number of knowledge elements it encompasses.

Building upon these per-query, per-element scores, the aggregate score of model  $M_t$  for a specific knowledge element  $k$  across the entire indexing dataset  $\mathcal{D}$  is computed as:

$$S_{\beta}^{\alpha}(M_t, \mathcal{D}, k) = \frac{1}{|\mathcal{D}^k|} \sum_{i=1}^N S_{\beta}^{\alpha}(M_t, \mathbf{x}_i, k)$$

where  $\mathcal{D}^k = \{(\mathbf{x}_i, y_i) \mid k \in \mathcal{K}_i\}$  denotes the subset of query-response pairs in which knowledge  $k$  is present in the knowledge set. A similar methodology is employed for the computation of capability scores.

### 3.4 Routing when inference

For an unseen query  $\mathbf{x}$  with its knowledge and capability sets, we compute the knowledge score  $KS$  and capability score  $CS$  for each candidate model  $M_t$  to guide routing. The knowledge score is given by:

$$KS^{\alpha}(M_t, \mathbf{x}) = \sum_{i=1}^{|\mathcal{K}_{\mathbf{x}}|} S_{\beta}^{\alpha}(M_t, \mathcal{D}, k_i) \frac{\alpha^{i-1}}{\sum_{m=1}^{|\mathcal{K}_{\mathbf{x}}|} \alpha^{m-1}}, \quad (1)$$

The capability score,  $CS^{\alpha}(M_t, \mathbf{x})$ , is computed analogously. Normalization across both knowledge and capability score calculations ensures that these two distinct types of scores are on a comparable scale, facilitating a balanced routing decision.

The final routing decision is determined by the following algorithm:

$$\mathcal{R}_{\mathcal{M}_T}(\mathbf{x}) = \arg \max_{M_t \in \mathcal{M}_T} (\gamma KS^{\alpha}(M_t, \mathbf{x}) + \delta CS^{\alpha}(M_t, \mathbf{x})) \quad (2)$$

which aims to identify the model with the highest weighted average of the knowledge and capability scores. A key advantage of this framework is its adaptability. New LLMs are efficiently integrated by evaluating them on  $\mathcal{D}$  to quantify their knowledge and capability scores, which are then used in routing. Similarly, when queries introduce novel knowledge, the LLMs' scores for this new knowledge can be computed and integrated, refining subsequent routing decisions.

## 4 Experiment

### 4.1 Dataset

In this section, we introduce our comprehensive dataset: **RouteMix**, which consist of the Index Set and Evaluation Set.

#### 4.1.1 Index Set

The term "Index Set" designates the dataset utilized during the development of our routing algorithm. Given that our methodology is parameter-free, this nomenclature serves to differentiate it from datasets conventionally used in training-dependent methods. The "Index Set" is thus employed primarily for characterizing and indexing the capabilities and knowledge of LLMs. To construct a sufficiently diverse "Index Set" for robust LLM profiling, we have curated 20 distinct datasets. These datasets span a wide array of domains and are instrumental in quantifying the specific knowledge and capabilities of each model. Comprehensive details regarding the statistics, data processing methodologies, and evaluation metrics for each dataset are presented in [Appx. §A](#).

#### 4.1.2 Evaluation Set

We incorporate four benchmarks that comprehensively evaluate the LLM as the evaluation set of

Method	MMLU-Pro	GPQA	BigGenBench	LiveBench	Avg.
<b>Single Large Language Model</b>					
Gemini-1.5-Pro	<b>82.83</b>	75.76	80.92	53.79	73.33
GPT-4o	79.71	74.24	<b>85.36</b>	49.62	72.23
Grok-2	80.14	76.26	83.66	53.26	73.33
Qwen2.5-Max	75.86	71.21	82.48	52.77	70.58
GLM-4-Plus	79.06	75.76	83.27	47.32	71.35
Nova-Pro	77.49	70.20	83.01	44.38	68.77
Llama-3.3-70B-Instruct	76.27	69.70	78.17	50.67	68.70
Qwen-2.5-72B-Instruct	75.41	73.23	82.61	49.83	70.27
Random	78.26	72.22	82.61	48.83	70.48
<b>Routing Algorithm</b>					
RouterDC	77.34	73.74	82.88	49.21	70.79
EmbedLLM	78.95	76.26	83.01	51.46	72.42
Routing by <i>Knowledge</i>	<u>80.99</u>	<b>78.28</b>	82.61	53.17	<u>73.76</u>
Routing by <i>Capability</i>	80.09	76.26	84.18	53.65	73.55
<i>Inference Dynamics</i>	80.85	<u>77.78</u>	<u>84.31</u>	<b>55.57</b>	<b>74.55</b>

Table 1: LLM routing results across four benchmarks are presented. The metrics we used are introduced in §4.2. The best performances are **bold-faced**, while the second-best performances are underlined. "Routing by Knowledge" denotes routing decisions made solely based on the knowledge score, whereas "Routing by Capability" refers to routing based only on the capability score. "Mixed Routing" indicates a simultaneous consideration of both scores during the routing process.

**RouteMix:** (i) MMLU-Pro (Wang et al., 2024b) spans 14 diverse domains and includes approximately 12,000 instances. (ii) GPQA (Rein et al., 2023) consists of multiple choice questions at the graduate level in subdomains of physics, chemistry, and biology. For our evaluation, we utilize the Diamond subset. (iii) BigGenBench (Kim et al., 2024) comprises 77 distinct tasks evaluating core abilities of LLM, with a total of 765 human-written instances. (iv) LiveBench (White et al., 2025) is a real-time updated benchmark with 18 tasks across 6 categories, including math, reasoning, coding, data analysis, language and instruction following. In the evaluation, we utilize the snapshot released on 2024-11-25.

## 4.2 Experiment Setup

For the candidate models, we select eight high-performing LLMs: Gemini-1.5-Pro (Reid et al., 2024), GPT-4o (Hurst et al., 2024), Grok-2, Qwen2.5-Max (Yang et al., 2024), GLM-4-Plus (Zeng et al., 2024), Nova-Pro (Intelligence, 2024), Llama-3.3-70B-Instruct (AI@Meta, 2024), and Qwen-2.5-72B-Instruct (Yang et al., 2024). To ensure a fair comparison when testing these models, all parameters and the input prompt are kept consistent across evaluations. To derive the Knowledge and Capability attributes, we employ GPT-4o-mini to generate these characteristics, and ablation study of auxiliary models is demonstrated in Appx. §C. Since generated attributes may include semantically similar phrases, we utilize MiniLM-

L6 (Wang et al., 2020) to consolidate Knowledge entries with a cosine similarity score greater than 0.6. Additionally, attributes with a frequency lower than 10 are filtered out and designated as "Other" entry. When the system encounters a query containing previously unseen attributes, these are also classified as "Other" entry. By default, for unconstrained routing, the parameters  $\alpha$  and  $\beta$  are set to 0.5 and 0, respectively. The weights for the Knowledge and Capability scores are both set to 1.0 by default. In terms of evaluation, the exact match score is employed for both the MMLU-Pro and GPQA datasets. For BigGenBench, we follow the methodology proposed by Sprague et al. (2025), using GPT-4o-mini as a language model-based judge. Instances receiving a score greater than 4 are classified as correct. For LiveBench, we adhere to the original evaluation script, and the metric is average score across six categories.

## 4.3 Capability and Knowledge Quantification

The performance of the candidate models on the Index Set is presented in Fig. 2. Generally, these models do not exhibit substantial performance distinctions when evaluated across the entire Index Set. However, their relative strengths become apparent on specific subsets, where different models tend to outperform one another. This observation suggests that the model pool consists of LLMs with broadly comparable overall abilities, yet with varying specializations.

Subsequent to the computation of average perfor-

385 mance scores, the top four models are selected for  
 386 more detailed analysis. Their respective capability  
 387 and knowledge scores are visualized in Fig. 1. For  
 388 clarity and simplification in this visualization, we  
 389 focus on the most frequently occurring knowledge  
 390 elements and capabilities within the Index Set. The  
 391 fact that the highest-scoring model changes with  
 392 the specific knowledge or capability further sub-  
 393 stantiates the premise: LLMs, even those exhibit-  
 394 ing similar aggregate performance levels, possess  
 395 distinct areas of specialized expertise.

#### 396 4.4 Optimal Routing

397 The optimal routing results, presented in Tab. 1,  
 398 highlight the clear superiority of our proposed rou-  
 399 ting strategies. Among these, our *Mixed Routing*  
 400 strategy, which combines both *Knowledge* and *Ca-*  
 401 *capability* scores, achieves the highest average per-  
 402 formance, outperforming the best single model,  
 403 Gemini-1.5-Pro, by a margin of 1.22. It also outper-  
 404 forms two strong routing baselines, demonstrating  
 405 the robust adaptability of our approach across the  
 406 comprehensive dataset. This strategy secures top  
 407 results on LiveBench and ranks second on GPQA  
 408 and BigGenBench, demonstrating the effectiveness  
 409 and versatility of our comprehensive routing algo-  
 410 rithm. Additionally, the Routing by *Knowledge*  
 411 and Routing by *Capability* approaches also deliver  
 412 strong results, consistently surpassing the best sin-  
 413 gle model and significantly outperforming random  
 414 routing on average. Notably, Routing by *Knowl-*  
 415 *edge* excels in knowledge-intensive tasks, achiev-  
 416 ing the best score on GPQA and the second-best  
 417 on MMLU-Pro. This underscores its ability to ef-  
 418 fectively direct queries requiring accurate factual  
 419 recall and nuanced domain understanding. Simi-  
 420 larly, Routing by *Capability* performs exceptionally  
 421 well on capability-driven benchmarks, particularly  
 422 on BigGenBench, highlighting the importance of  
 423 leveraging a model’s inherent strengths in complex  
 424 reasoning and generation tasks. Both approaches  
 425 play an integral role in the success of the *Mixed*  
 426 *Routing* system.

427 These findings also emphasize that no single  
 428 LLM universally dominates across all tasks. Mod-  
 429 els like Gemini-1.5-Pro and GPT-4o exhibit vary-  
 430 ing strengths, further validating the necessity and  
 431 advantages of intelligent LLM routing systems.

#### 432 4.5 Routing with Constraints

433 To investigate the system’s performance under vary-  
 434 ing cost constraints, we systematically adjusted the

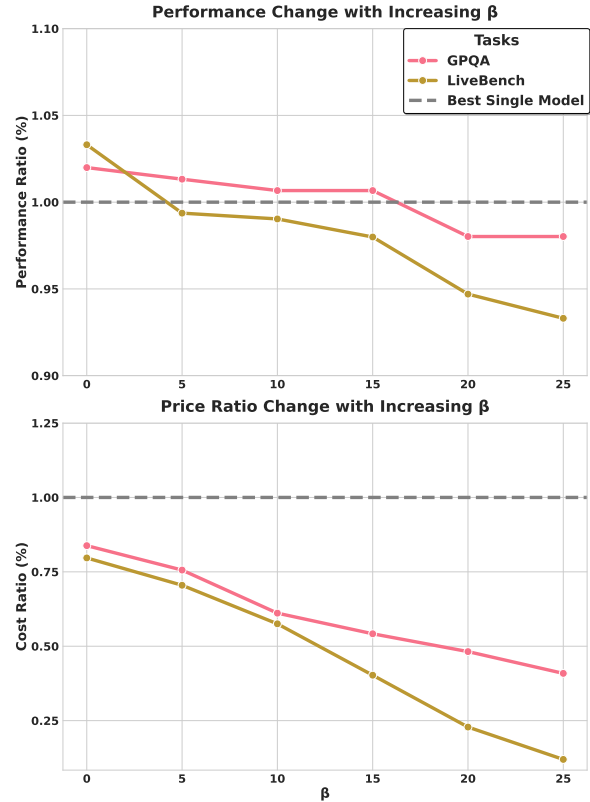


Figure 3: Performance Ratio (%) and Cost Ratio (%) variation on GPQA and LiveBench. The "Best Single Model" refers to the most performant LLM for each task.

435  $\beta$  parameter, maintaining all other experimental  
 436 configurations as previously defined. The evaluation  
 437 employed two distinct metrics. The first metric,  
 438 termed **Performance Ratio**, quantifies the efficacy  
 439 of the *Mixed Routing* strategy. This is calculated  
 440 as the ratio of the performance achieved by *Mixed*  
 441 *Routing* to that of the best-performing single can-  
 442 didate LLM on the respective benchmark. The  
 443 second metric, **Cost Ratio**, assesses the economic  
 444 efficiency of the routing algorithm. It is defined as  
 445 the total cost incurred by the routing process (en-  
 446 compassing both knowledge generation and capa-  
 447 bility assessment costs) relative to the operational  
 448 cost of the best-performing single LLM.

449 The empirical results of this sensitivity analysis  
 450 are depicted in Fig. 3. In scenarios without strin-  
 451 gent price constraints (i.e.,  $\beta = 0$ ), our routing  
 452 system demonstrates superior performance com-  
 453 pared to the best single model, while operating at  
 454 approximately 80% of the latter’s budget. As the  
 455  $\beta$  parameter is incrementally increased, thereby  
 456 prioritizing cost reduction, the operational cost of  
 457 the routing algorithm decreases significantly. Con-  
 458 currently, the system maintains a competitive per-

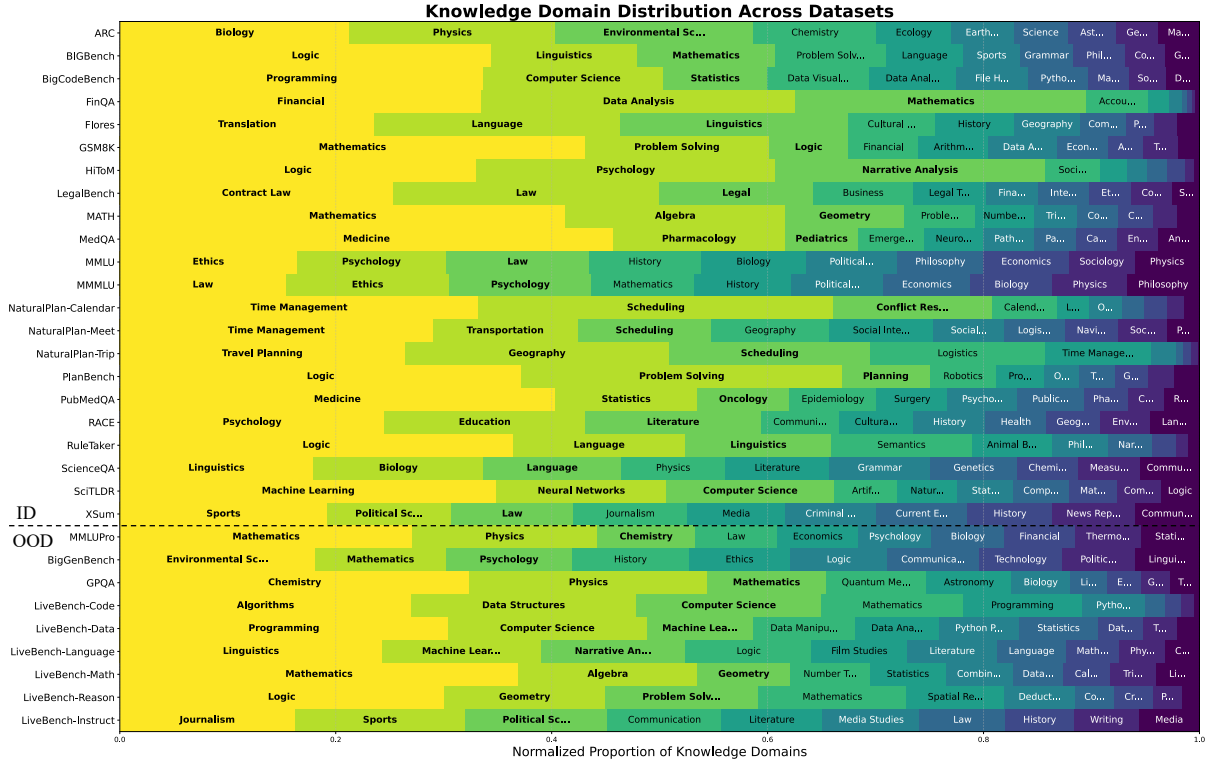


Figure 4: Distribution of knowledge domains across 24 datasets in **RouteMix**. The In-Domain (ID) subset is utilized for quantifying *Knowledge* and *Capability*, while the Out-of-Domain (OOD) subset is employed for evaluating the routing algorithm. Dataset labels including "LiveBench" indicate subsets of the LiveBench benchmark, and labels including "NaturalPlan" similarly denote subsets of the NaturalPlan benchmark. The algorithm to compute the normalized proportion is included in [Appx. §B](#).

459 formance level relative to the best single model. 481  
 460 Notably, at a  $\beta$  value of 15, our routing algorithm 482  
 461 achieves performance nearly equivalent to the best 483  
 462 single model but utilizes only approximately half 484  
 463 the associated cost. 485

464 An interesting observation is the differential sensi- 486  
 465 tivity of benchmarks to changes in  $\beta$ . Specifically, 487  
 466 the performance and cost metrics for LiveBench, 488  
 467 a text generation benchmark, exhibit more pro- 489  
 468 nounced variations in response to adjustments in  $\beta$  490  
 469 compared to those observed for GPQA, a question- 491  
 470 answering benchmark. This suggests that text gen- 492  
 471 eration tasks are more sensitive to the price penalty 493  
 472 imposed by  $\beta$  than QA tasks. 494

## 473 5 Analysis 495

### 474 5.1 Model Selection 496

475 The distribution of model selections under vari- 497  
 476 ous conditions is illustrated in [Fig. 5](#). Consistent 498  
 477 with findings in previous works ([Chen et al., 2024](#); 499  
 478 [Frick et al., 2025](#)), cost-efficient models are in- 500  
 479 frequently selected in optimal routing scenarios; 501  
 480 instead, the strategy predominantly converges to- 502  
 503  
 504

wards higher-performing models. For comprehen- 500  
 sive benchmarks such as BigGenBench, our ap- 501  
 proach primarily routes queries to expensive yet 502  
 high-performing models like GPT-4o and Grok-2, 503  
 reflecting a tendency to leverage top-tier capabili- 504  
 ties for broad-ranging tasks. Conversely, for task 505  
 sets demanding highly specialized capabilities, the 506  
 routing algorithm typically assigns queries directly 507  
 to the most proficient model. For instance, within 508  
 the coding subset of LiveBench, 91% of queries 509  
 are routed to Qwen-Max, which demonstrates the 510  
 strongest coding capabilities. This model’s leading 511  
 performance in coding is further corroborated by 512  
 its results on BigCodeBench and its specific Cod- 513  
 ing capability score, as detailed in [Fig. 1](#) and [Fig. 2](#), 514  
 respectively. These observations collectively indi- 515  
 cate that our routing algorithm effectively directs 516  
 queries to the most suitable models based on spe- 517  
 cific task demands. 518

In the context of cost-constrained routing, an 500  
 increasing cost penalty prompts the router to pro- 501  
 gressively shift its selections from expensive, top- 502  
 performing models towards more affordable, albeit 503  
 less powerful, alternatives. 504

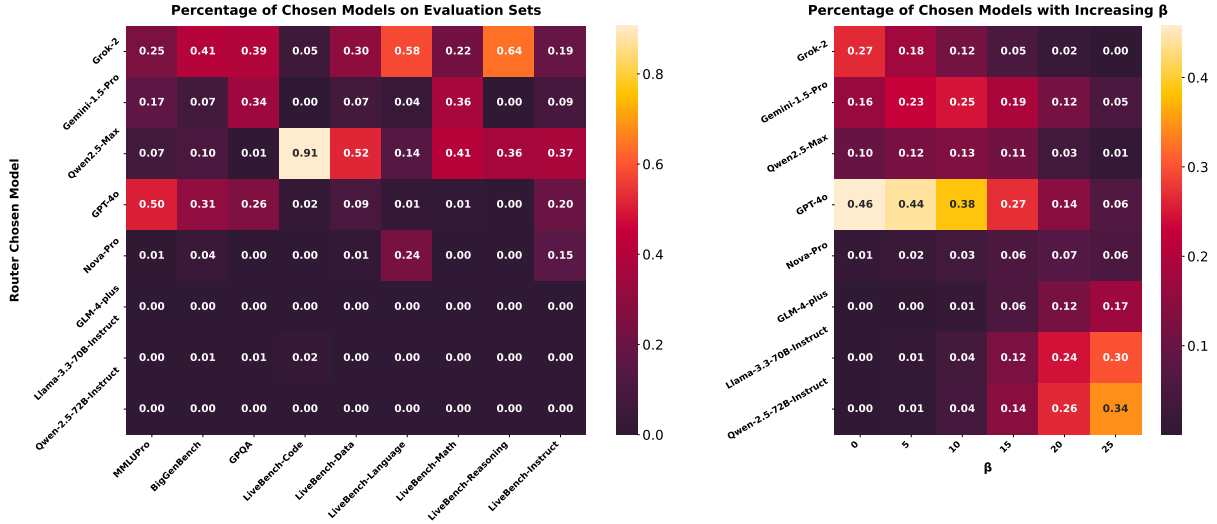


Figure 5: Comparative distribution of router-selected models. Lighter colors signify a higher selection ratio for a given model. The left panel details model selection across evaluation benchmarks using the Optimal *Mixed Routing* strategy. The right panel illustrates the impact of an increasing cost penalty coefficient ( $\beta$ ) on the model selection distribution.

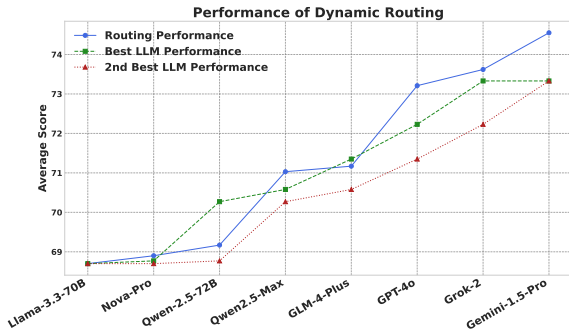


Figure 6: Routing Performance (%) in Dynamic LLM Pools.

## 5.2 Knowledge Distribution

As shown in Fig. 4, the distribution of generated knowledge highlights the *RouteMix* benchmark’s comprehensive span of knowledge domains, ranging from highly specific academic areas to practical applications. On datasets with broad knowledge requirements, such as MMLU-Pro, the generated knowledge exhibits a relatively balanced distribution. For benchmarks targeting one or two specific domains, like MATH-500, the model typically generates more fine-grained knowledge components related to the core domain. This facilitates a more nuanced quantification of the model’s domain-specific knowledge.

## 5.3 Dynamics Routing

In this section, we investigate the scalability of our framework with respect to dynamic LLM pools.

The corresponding results are presented in Fig. 6. The x-axis in this figure represents the progressive addition of specific new models to the LLM candidate pool. Initially, the pool consists solely of Llama-3.3-70B; subsequently, one new model is added to the candidate pool at each increment along the x-axis. Notably, our routing algorithm consistently maintains a top-2 performance ranking and surpasses the best single model across the five evaluated candidate pool configurations. This outcome demonstrates the robust scalability of our framework when new models are introduced, crucially without the need for any additional training.

## 6 Conclusions

This paper introduces **InferenceDynamics**, a scalable and adaptable LLM routing framework that quantifies model capabilities and domain-specific knowledge to match queries with the most suitable LLMs. Evaluated on the new comprehensive *RouteMix* benchmark, *InferenceDynamics* demonstrated superior performance, outperforming the best single LLM by 1.22 on average and achieving comparable results at approximately half the cost under budget constraints. Key contributions include the **RouteMix** dataset for evaluating generalization and the **InferenceDynamics** algorithm, which generalizes to unseen queries and effectively routes them within dynamic model pools without retraining. Our work enables more efficient and tailored utilization of the diverse LLM ecosystem.

## 552 Limitations

553 Despite the promising results and the robust design  
554 of InferenceDynamics, several limitations warrant  
555 discussion and offer avenues for future research:

556 **Niche Suitability for Highly Constrained En-**  
557 **vironments** InferenceDynamics is engineered  
558 for scalability and adaptability, demonstrating its  
559 strengths when dealing with a large, diverse, and  
560 evolving pool of LLMs, or when new capability and  
561 knowledge domains are frequently encountered.  
562 However, in scenarios characterized by a very lim-  
563 ited and static set of LLMs and a narrowly de-  
564 fined, unchanging task scope, a dedicated learning-  
565 based routing approach (e.g., a fine-tuned classifier)  
566 might be more appropriate or yield marginally su-  
567 perior, hyper-specialized performance. Our frame-  
568 work prioritizes generalizability and efficient adap-  
569 tation to dynamic conditions, which is a differ-  
570 ent niche than hyper-optimization for small, fixed-  
571 scope problems.

572 **Benchmark-Driven Evaluation vs. Real-World**  
573 **Application Complexity** The current evaluation  
574 of InferenceDynamics relies on the comprehensive  
575 RouteMix dataset, which is composed of various  
576 established benchmarks. While these benchmarks  
577 cover a wide array of tasks and domains, they may  
578 not fully capture the intricacies and dynamic nature  
579 of real-world application systems. For instance,  
580 the utility and performance of InferenceDynamics  
581 in more complex, interactive systems like multi-  
582 agent environments, where task allocation might  
583 depend on evolving collaborative states, have not  
584 been explicitly tested. Exploring the deployment  
585 and effectiveness of InferenceDynamics in such  
586 real-application scenarios remains an important di-  
587 rection for future work.

588 Addressing these limitations will be crucial for  
589 broadening the applicability and enhancing the ro-  
590 bustness of InferenceDynamics and similar LLM  
591 routing frameworks.

## 592 Ethics Statement

593 Our study utilizes publicly available datasets and  
594 accesses Large Language Models (LLMs) through  
595 their respective APIs. The ethical considerations  
596 pertaining to this research are as follows:

597 **Datasets:** This research exclusively employs pub-  
598 licly available datasets, strictly for academic re-  
599 search purposes. We affirm that no personally iden-  
600 tifiable information or private data was involved in

our study.

**LLM APIs:** Our application of LLMs via APIs  
rigorously conforms to the policies set forth by  
the API providers. This includes adherence to fair  
use guidelines and respect for intellectual property  
rights.

**Transparency:** In line with standard academic  
research practices, we provide detailed descriptions  
of our methodology and the prompts utilized in our  
experiments. Furthermore, the source code for this  
research will be made publicly available upon the  
acceptance of this paper.

## References

- Pranjal Aggarwal, Aman Madaan, Ankit Anand, Sriv-  
idya Pranavi Potharaju, Swaroop Mishra, Pei Zhou,  
Aditya Gupta, Dheeraj Rajagopal, Karthik Kappa-  
ganthu, Yiming Yang, Shyam Upadhyay, Manaal  
Faruqui, and Mausam. 2024. [Automix: Automati-  
cally mixing language models](#). In *Advances in Neu-  
ral Information Processing Systems 38: Annual Con-  
ference on Neural Information Processing Systems  
2024, NeurIPS 2024, Vancouver, BC, Canada, De-  
cember 10 - 15, 2024*.
- AI@Meta. 2024. [Llama 3 model card](#).
- Isabel Cachola, Kyle Lo, Arman Cohan, and Daniel S.  
Weld. 2020. [TLDR: extreme summarization of sci-  
entific documents](#). In *Findings of the Association for  
Computational Linguistics: EMNLP 2020, Online  
Event, 16-20 November 2020*, volume EMNLP 2020  
of *Findings of ACL*, pages 4766–4777. Association  
for Computational Linguistics.
- Lingjiao Chen, Matei Zaharia, and James Zou. 2023.  
[Frugalpt: How to use large language models while  
reducing cost and improving performance](#). *CoRR*,  
abs/2305.05176.
- Shuhao Chen, Weisen Jiang, Baijiong Lin, James T.  
Kwok, and Yu Zhang. 2024. [Routerdc: Query-based  
router by dual contrastive learning for assembling  
large language models](#). In *Advances in Neural In-  
formation Processing Systems 38: Annual Confer-  
ence on Neural Information Processing Systems 2024,  
NeurIPS 2024, Vancouver, BC, Canada, December  
10 - 15, 2024*.
- Zhijun Chen, Jingzheng Li, Pengpeng Chen, Zhuoran Li,  
Kai Sun, Yuankai Luo, Qianren Mao, Dingqi Yang,  
Hailong Sun, and Philip S. Yu. 2025. [Harnessing  
multiple large language models: A survey on LLM  
ensemble](#). *CoRR*, abs/2502.18036.
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena  
Shah, Iana Borova, Dylan Langdon, Reema Moussa,  
Matt Beane, Ting-Hao Kenneth Huang, Bryan R.  
Routledge, and William Yang Wang. 2021. [Finqa:  
A dataset of numerical reasoning over financial data](#).

654	In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021</i> , pages 3697–3711. Association for Computational Linguistics.		
655			
656			
657			
658			
659	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. <i>arXiv:1803.05457v1</i> .		
660			
661			
662			
663			
664	Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. In <i>Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020</i> , pages 3882–3890. ijcai.org.		
665			
666			
667			
668			
669	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. <i>CoRR</i> , abs/2110.14168.		
670			
671			
672			
673			
674			
675	Xiangxiang Dai, Jin Li, Xutong Liu, Anqi Yu, and John C. S. Lui. 2024. Cost-effective online multi-llm selection with versatile reward models. <i>CoRR</i> , abs/2405.16587.		
676			
677			
678			
679	Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Rühle, Laks V. S. Lakshmanan, and Ahmed Hassan Awadallah. 2024. Hybrid LLM: cost-efficient and quality-aware query routing. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.		
680			
681			
682			
683			
684			
685			
686	Evan Frick, Connor Chen, Joseph Tennyson, Tianle Li, Wei-Lin Chiang, Anastasios N. Angelopoulos, and Ion Stoica. 2025. Prompt-to-leaderboard. <i>CoRR</i> , abs/2502.14855.		
687			
688			
689			
690	Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The flores-101 evaluation benchmark for low-resource and multilingual machine translation. <i>Trans. Assoc. Comput. Linguistics</i> , 10:522–538.		
691			
692			
693			
694			
695			
696			
697	Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, K. Aditya, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John J. Nay, Jonathan H. Choi, Kevin Tobia, Margaret Hagan, Megan Ma, Michael A. Livermore, Nikon Rasumov-Rahe, Nils Holzenberger, Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny Gandhi, Tom Zur, Varun Iyer, and Zehua Li. 2023. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language		
698			
699			
700			
701			
702			
703			
704			
705			
706			
707			
708			
709			
710			
		models. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	711
			712
			713
			714
		Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. <i>CoRR</i> , abs/2009.03300.	715
			716
			717
			718
		Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	719
			720
			721
			722
			723
			724
			725
		Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll L. Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, and Dane Sherburn. 2024. Gpt-4o system card. <i>CoRR</i> , abs/2410.21276.	726
			727
			728
			729
			730
			731
			732
			733
			734
			735
			736
			737
			738
			739
			740
			741
			742
			743
			744
			745
			746
			747
			748
			749
			750
			751
			752
			753
			754
			755
			756
			757
		Amazon Artificial General Intelligence. 2024. The amazon nova family of models: Technical report and model card. <i>Amazon Technical Reports</i> .	758
			759
			760
		Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 14165–14178. Association for Computational Linguistics.	761
			762
			763
			764
			765
			766
			767
			768

769	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. <a href="#">What disease does this patient have? A large-scale open domain question answering dataset from medical exams</a> . <i>CoRR</i> , abs/2009.13081.	
770		
771		
772		
773		
774	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W. Cohen, and Xinghua Lu. 2019. <a href="#">Pubmedqa: A dataset for biomedical research question answering</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 2567–2577. Association for Computational Linguistics.	
775		
776		
777		
778		
779		
780		
781		
782		
783	Seungone Kim, Juyoung Suk, Ji Yong Cho, Shayne Longpre, Chaeun Kim, Dongkeun Yoon, Guijin Son, Yejin Choi, Sheikh Shafayat, Jinheon Baek, Sue Hyun Park, Hyeonbin Hwang, Jinkyung Jo, Hyowon Cho, Haebin Shin, Seongyun Lee, Hanseok Oh, Noah Lee, Namgyu Ho, Se June Joo, Miyoung Ko, Yoonjoo Lee, Hyungjoo Chae, Jamin Shin, Joel Jang, Seonghyeon Ye, Bill Yuchen Lin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. <a href="#">The biggen bench: A principled benchmark for fine-grained evaluation of language models with language models</a> . <i>CoRR</i> , abs/2406.05761.	
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796	Steven Kolawole, Don Kurian Dennis, Ameet Talwalkar, and Virginia Smith. 2024. <a href="#">Revisiting cascaded ensembles for efficient inference</a> . <i>CoRR</i> , abs/2407.02348.	
797		
798		
799		
800	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard H. Hovy. 2017. <a href="#">RACE: large-scale reading comprehension dataset from examinations</a> . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017</i> , pages 785–794. Association for Computational Linguistics.	
801		
802		
803		
804		
805		
806		
807		
808	Jiawei Li, Yizhe Yang, Yu Bai, Xiaofeng Zhou, Yinghao Li, Huashan Sun, Yuhang Liu, Xingpeng Si, Yuhao Ye, Yixiao Wu, Yiguan Lin, Bin Xu, Ren Bowen, Chong Feng, Yang Gao, and Heyan Huang. 2024a. <a href="#">Fundamental capabilities of large language models and their applications in domain scenarios: A survey</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 11116–11141. Association for Computational Linguistics.	
809		
810		
811		
812		
813		
814		
815		
816		
817		
818		
819	Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. 2024b. <a href="#">More agents is all you need</a> . <i>CoRR</i> , abs/2402.05120.	
820		
821		
822	Moxin Li, Yong Zhao, Yang Deng, Wenxuan Zhang, Shuaiyi Li, Wenya Xie, See-Kiong Ng, and Tat-Seng Chua. 2024c. <a href="#">Knowledge boundary of large language models: A survey</a> . <i>CoRR</i> , abs/2412.12472.	
823		
824		
825		
	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. <a href="#">Learn to explain: Multimodal reasoning via thought chains for science question answering</a> . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	826
		827
		828
		829
		830
		831
		832
		833
		834
	Andrea Matarazzo and Riccardo Torlone. 2025. <a href="#">A survey on large language models with some insights on their capabilities and limitations</a> . <i>CoRR</i> , abs/2501.04040.	835
		836
		837
		838
	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. <a href="#">Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018</i> , pages 1797–1807. Association for Computational Linguistics.	839
		840
		841
		842
		843
		844
		845
		846
	Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E. Gonzalez, M. Waleed Kadous, and Ion Stoica. 2025. <a href="#">Routellm: Learning to route llms from preference data</a> . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	847
		848
		849
		850
		851
		852
		853
	OpenAI. 2025. <a href="#">Introducing gpt-5</a> . <a href="https://openai.com/index/introducing-gpt-5/">https://openai.com/index/introducing-gpt-5/</a> . Blog post.	854
		855
		856
	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Abbas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal, Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren Sezener, and et al. 2024. <a href="#">Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context</a> . <i>CoRR</i> , abs/2403.05530.	857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. <a href="#">GPQA: A graduate-level google-proof q&amp;a benchmark</a> . <i>CoRR</i> , abs/2311.12022.	877
		878
		879
		880
		881
	Zayne Rea Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann	882
		883

884	Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. 2025. <a href="#">To cot or not to cot? chain-of-thought helps mainly on math and symbolic reasoning</a> . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	942
885		943
886		944
887		945
888		946
889		947
890	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. 2023. <a href="#">Challenging big-bench tasks and whether chain-of-thought can solve them</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13003–13051. Association for Computational Linguistics.	948
891		949
892		950
893		951
894		952
895		953
896		954
897		955
898		
899	Karthik Valmeekam, Matthew Marquez, Alberto Olmo Hernandez, Sarath Sreedharan, and Subbarao Kambhampati. 2023. <a href="#">Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change</a> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	956
900		957
901		958
902		959
903		960
904		
905		
906		
907		
908	Qineng Wang, Zihao Wang, Ying Su, Hanghang Tong, and Yangqiu Song. 2024a. <a href="#">Rethinking the bounds of LLM reasoning: Are multi-agent discussions the key?</a> In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 6106–6131. Association for Computational Linguistics.	961
909		962
910		963
911		964
912		965
913		966
914		967
915		968
916		969
917		
918		
919		
920		
921		
922		
923	Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. <a href="#">Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers</a> . In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	970
924		971
925		972
926		973
927		974
928		975
929		976
930		977
931		978
932		979
933		980
934		981
935		
936		
937		
938	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024b. <a href="#">Mmlu-pro: A more robust and challenging multi-task language understanding benchmark</a> . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	982
939		983
940		984
941		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

1001 with multi-objective optimal consideration. *CoRR*,  
1002 abs/2410.08014.

1003 Yi-Kai Zhang, Ting-Ji Huang, Yao-Xiang Ding, De-  
1004 Chuan Zhan, and Han-Jia Ye. 2023. [Model spider: Learning to rank pre-trained models efficiently](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

1010 Yi-Kai Zhang, De-Chuan Zhan, and Han-Jia Ye. 2025.  
1011 [Capability instruction tuning: A new paradigm for dynamic llm routing](#). *Preprint*, arXiv:2502.17282.

1013 Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang,  
1014 Xinyun Chen, Minmin Chen, Azade Nova, Le Hou,  
1015 Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, and  
1016 Denny Zhou. 2024. [NATURAL PLAN: benchmarking llms on natural language planning](#). *CoRR*,  
1017 abs/2406.04520.

1019 Richard Zhuang, Tianhao Wu, Zhaojin Wen, Andrew Li,  
1020 Jiantao Jiao, and Kannan Ramchandran. 2025. [Embedlm: Learning compact representations of large language models](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.

1025 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu,  
1026 Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani  
1027 Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al.  
1028 2024. [Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions](#). *arXiv preprint arXiv:2406.15877*.

## A Benchmark Overview Table

Table 2: Overview of Benchmarks, Data Processing, Prompts, and Metrics

Benchmark Name	Data Processing Manner	Prompt Type	Metric Used
ARC (Clark et al., 2018)	Sample 500 instances according to the portion of ARC-Easy and ARC-Challenge.	Zero-shot DA	Accuracy
BigBench-Hard (Suzgun et al., 2023)	Sample 40 instances from each category except <i>web_of_lies</i> , to avoid collision with LiveBench. Formulate into MCQA for Yes/No and QA question. Remain the free-response question unchanged.	Zero-shot CoT	Exact Match (EM)
BigCodeBench (Zhuo et al., 2024)	We directly use the BigCodeBench-Hard subset, with 148 instances.	DA for code completion	Pass@1
FinQA (Chen et al., 2021)	Sample 500 instances from the dataset.	CoT from Sprague et al. (2025)	Exact Match(EM)
Flores200 (Goyal et al., 2022)	We incorporate the top10 commonly used language except for English. And sample 100 instances for each language.	Translation Prompt	Chrf++ (Goyal et al., 2022)
GSM8K (Cobbe et al., 2021)	Sample 500 instances from the dataset.	CoT from Sprague et al. (2025)	Exact Match(EM).
HiToM (Wu et al., 2023)	Sample 500 instances under CoT settings.	CoT from Official Repo	Accuracy
LegalBench (Guha et al., 2023)	Sample 4 instances from each category except for short answering task, resulting in 616 instances.	Few-shot DA	Accuracy
MATH (Hendrycks et al., 2021)	We use the subset MATH-500.	CoT from Sprague et al. (2025)	Exact Match(EM)
MedQA (Jin et al., 2020)	Sample 500 instances from the dataset	DA	Accuracy
MMLU (Hendrycks et al., 2020)	We sample instances according to the portion of different categories, and make sure each category has at least 10 instances. Resulting in 1262 instances.	DA	Accuracy
MMMLU (Hendrycks et al., 2020)	We sample 100 instances for all languages except for English. Result in 1400 instances.	DA	Accuracy
NaturalPlan (Zheng et al., 2024)	Sample 200 instances from each subset, including scheduling, calendar meeting, and trip planning.	DA	Accuracy
PlanBench (Valmeekam et al., 2023)	Use the subset of PlanGeneration in BlocksWorld.	DA	Accuracy

Continued on next page...

Table 2 – continued from previous page

Benchmark Name	Data Processing Manner	Prompt Type	Metric Used
PubMedQA (Jin et al., 2019)	Sample 500 instances from original dataset	DA	Accuracy
RACE (Lai et al., 2017)	Sample 500 instances from original dataset	DA	Accuracy
RuleTaker (Clark et al., 2020)	Sample 500 instances from original dataset	DA	Accuracy
ScienceQA (Lu et al., 2022)	Sample 500 instances which don't have corresponding picture.	DA	Accuracy
SciTLDR (Cachola et al., 2020)	Directly use the test set	Summarization Prompt	RogueL.
XSum (Narayan et al., 2018)	Sample 500 instances for the dataset	Summarization Prompt	RogueL

Specifically, when quantifying the capability and knowledge of LLMs for translation and summarization tasks, we establish a performance threshold. An output is considered correct if its evaluation score or relevant metric exceeds this threshold.

## B Knowledge Domain Distribution

The dataset's knowledge domain distribution is determined by a weighted rank approach. For each domain  $D \in \mathcal{D}$  (where  $\mathcal{D}$  is the set of all unique domains), its frequency at each rank  $r$  (denoted  $F_{D,r}$ , for  $r = 1, \dots, N$ ) is multiplied by a corresponding rank weight  $W_r$  (typically  $W_r = 1/r$ ). These products are summed to yield a weighted score  $S_D$ :

$$S_D = \sum_{r=1}^N (F_{D,r} \times W_r)$$

The final distribution percentage  $P_D$  for each domain is then its  $S_D$  normalized by the sum of all domain weighted scores ( $S_{\text{total}} = \sum_{D' \in \mathcal{D}} S_{D'}$ ), expressed as a percentage:

$$P_D = \left( \frac{S_D}{\sum_{D' \in \mathcal{D}} S_{D'}} \right) \times 100\%$$

This method ensures higher-ranked domain occurrences contribute more significantly, with all  $P_D$  summing to 100%.

## C Robustness to Auxiliary Model Choice

To test whether our framework's performance is overly dependent on a single auxiliary model, we performed routing using attributes generated by three different models: Qwen-2.5-7b-instruct, Gemma-3-12b-it, and GPT-4o-mini. The key finding is that while there are minor performance variations, using any of the tested auxiliary models—from the lightweight Qwen-2.5-7b to the powerful GPT-4o-mini results in an average performance that surpasses the "Best Single Model" baseline. This demonstrates that our framework is not overly sensitive to the specific biases or minor accuracy differences of the auxiliary model, confirming its robustness. Regarding the concern about additional costs, we analyzed the cost of using each auxiliary model for attribute annotation relative to the total cost of routing (i.e., the cost of the final inference). As the data shows, the cost of attribute annotation is marginal, consistently accounting for less than 3% of the total routing expenditure, and often less than 1% when using efficient open-source models. This confirms that the cost of the auxiliary model is a negligible component of the overall operational budget.

Method	MMLUPro	GPQA	BigGenBench	LiveBench	Avg.
Best Single Model	82.83	75.76	80.92	53.79	73.33
Random	78.26	72.22	82.61	48.83	70.48
Qwen-2.5-7b-instruct	81.33	76.26	82.75	53.01	73.34
Gemma-3-12b-it	82.06	76.77	83.92	54.50	74.31
GPT-4o-mini	80.85	<b>77.78</b>	<b>84.31</b>	<b>55.57</b>	<b>74.55</b>

Table 3: InferenceDynamics results derived by different auxiliary models. The "Method" column indicates which model was used as the auxiliary model to generate the attributes for routing.

	MMLUPro	GPQA	BigGenBench	LiveBench
Qwen-2.5-7b-instruct	0.846%	0.623%	0.879%	1.044%
Gemma-3-12b-it	0.831%	0.742%	0.865%	0.846%
GPT-4o-mini	2.287%	1.834%	2.522%	2.335%

Table 4: Cost (%) of models to generate attributes across different benchmarks.

## D Further analysis over Model-SAT

While both methods use a predefined set of capabilities, InferenceDynamics offers several distinct advantages in terms of functionality, efficiency, and adaptability.

- 1. Flexible Cost-Performance Control.** The framework introduces a penalty parameter  $\beta$ , which enables explicit and adjustable trade-offs between inference cost and accuracy. This level of control is absent in Model-SAT and allows users to dynamically balance budget and performance requirements depending on task demands.
- 2. Superior Inference Efficiency.** InferenceDynamics significantly reduces computational overhead compared to Model-SAT’s routing mechanism. When the number of candidate models is  $\mathcal{M}$  and the routed model has  $\mathcal{N}$  parameters, Model-SAT requires  $\mathcal{M}$  forward passes through the router plus an additional inference step, yielding a total complexity of  $\mathcal{O}(\mathcal{M} + \mathcal{N})$ . InferenceDynamics only needs inference of the selected model, with complexity  $\mathcal{O}(\mathcal{N})$ , making the approach more scalable as the number of models increases.
- 3. Adaptability Without Retraining.** InferenceDynamics adapts to new domains without retraining. While Model-SAT must update its capability instructions, evaluate all LLMs on new data, and retrain its router, InferenceDynamics incorporates novel capabilities simply by expanding its index dataset with inference scores. This allows the system to evolve efficiently as tasks and domains shift, without additional retraining costs.

## E BigGenBench Evaluation

1077

Following Sprague et al. (2025), we employ GPT-4o-mini as LLM-as-a-Judge to evaluate the BigGenBench, and instances with a score larger than 4 is considered correct. The specific prompt is shown below:

1078

1079

1080

### Prompt for evaluation BigGenBench

#### Task Description:

An instruction (might include an Input inside it), a response to evaluate, a reference answer that gets a score of 5, and a score rubric representing a evaluation criteria are given.

1. Write a detailed feedback that assess the quality of the response strictly based on the given score rubric, not evaluating in general.
2. After writing a feedback, write a score that is an integer between 1 and 5. You should refer to the score rubric.
3. The output format should look as follows: "Feedback: (write a feedback for criteria) [RESULT] (an integer number between 1 and 5)"
4. Please do not generate any other opening, closing, and explanations.

The instruction to evaluation:  
example question

Response to evaluate:  
example solution

Reference Answer (Score 5):  
reference score

Score Rubrics:  
Criteria:  
criteria

Description of a Score 1 response:  
score1 description

Description of a Score 2 response:  
score2 description

Description of a Score 3 response:  
score3 description

Description of a Score 4 response:  
score4 description

Description of a Score 5 response:  
score5 description

Feedback:  
Remember, you must strictly evaluate the response based on the given score rubric, and finish your output in the format of "(...) [RESULT] <score>", where <score> is a number between 1 and 5.

1081

## F Prompt of Knowledge and Capability Generation

The specific prompt for knowledge and capability generation is shown below:

### Prompt for evaluation BigGenBench

The capabilities of Language Models include the following:

- Reasoning: Ability to logically analyze information, draw conclusions, and make inferences.
- Comprehension (Applicable to queries involving long passage comprehension): Understanding and interpreting the meaning, context, and nuances of extended or complex long-context text, such as lengthy documents, multi-paragraph inputs, or intricate narratives.
- Instruction Following (Applicable to queries involving several constraints): Accurately adhering to explicit user-provided guidelines, constraints, or formatting requirements specified within the query.
- Agentic: Capacity related to agent-like behavior, such as actively formulating plans, strategically deciding steps, and autonomously identifying solutions or actions to achieve specific goals or complex tasks.
- Knowledge Retrieval: Accessing and presenting accurate factual information from pre-existing knowledge.
- Coding: Generating, interpreting, or debugging computer programs and scripts.
- In-context Learning: Learning from examples or context provided within the current interaction without additional training.
- Multilingual (Must rank it in top3 when queries involving languages other than English): Understanding, generating, or translating content accurately across multiple languages.

Given the Query below:

1. Identify and list the \*LLM Capabilities\* from the definitions above that are directly and significantly required to effectively address the query.
  2. Identify and list the general \*Knowledge Domains\* (e.g., categories, subject areas) most pertinent to solving the problem presented in the query.
- List the selected Capabilities first, ranked from most important to least important. Then, list the identified Knowledge Domains, also ranked from most important to least important. \*Do not provide any justification or explanation\* for your selections or rankings.

Example:

Query: "Solve the following financial problem efficiently and clearly. Output the final answer as: boxedanswer. Where [answer] is just the final number or expression that solves the problem. Keep the answer to five decimal places if it is a number, and do not use percentages; keep the decimal format.

Problem: what is the net change in net revenue during 2016 for Entergy Mississippi, Inc.? the 2015 net revenue of amount (in millions) is 696.3; the 2016 net revenue of amount (in millions) is 705.4; Entergy Mississippi, Inc."

Capabilities: Reasoning, Knowledge retrieval

Knowledge: 1. Financial 2. Math 3. Data Analysis ... Query: input prompt