

# LEARNING COMPACT REPRESENTATIONS OF LLM ABILITIES VIA ITEM RESPONSE THEORY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent years have witnessed a surge in the number of large language models (LLMs), yet efficiently managing and utilizing these vast resources remains a significant challenge. In this work, we explore how to learn compact representations of LLM abilities that can facilitate downstream tasks, such as model routing and benchmark prediction. We frame this problem as estimating the probability that a given model will correctly answer a specific query. Inspired by the item response theory (IRT) in psychometrics, we model this probability as a function of three key factors: (i) the model’s multi-skill ability embedding  $\theta$ , (ii) the query’s discrimination vector  $\alpha$  that separates models of differing skills, and (iii) the query’s difficulty scalar  $\beta$ . To learn these parameters jointly, we introduce a Mixture-of-Experts (MoE) network that couples model- and query-level embeddings. Extensive experiments demonstrate that our approach leads to state-of-the-art performance in both model routing and benchmark accuracy prediction. Moreover, analysis validates that the learned parameters encode meaningful, interpretable information about model capabilities and query characteristics.

## 1 INTRODUCTION

Recent years have seen an explosion of large language models (LLMs), spanning from massive, general-purpose systems to small, task-specialized ones. As of August 2025, Hugging Face lists over 97,000 text-generation models with at least 6B parameters. This proliferation, however, introduces a significant challenge: how to *efficiently* manage and utilize such a vast, rapidly expanding ecosystem.

A crucial step toward addressing this challenge is to construct *compact* representations of models’ abilities (Zhuang et al., 2025). By encoding each model’s strengths and weaknesses, such representations can facilitate a range of downstream applications. For example, in *model routing*, they allow for assessing a model’s suitability for a given query, so that queries can be assigned to the most appropriate model within a candidate pool (Shnitzer et al., 2023; Lu et al., 2024; Jitkrittum et al., 2025). This not only helps balance performance and cost (Ong et al., 2024; Feng et al., 2024; Wang et al., 2025; Frick et al., 2025), but also empowers ensembles of smaller models to effectively compete with large proprietary systems (Zhang et al., 2025e; Pan et al., 2025). Another use case is *benchmark prediction* (Polo et al., 2024b; Zhang et al., 2025a). Traditionally, benchmarks are designed to compare the abilities of different models on specific tasks. However, the recent surge in their number<sup>1</sup> makes it impractical for exhaustive evaluation. Compact representations of model abilities offer an alternative, enabling efficient, scalable LLM evaluation.

In this work, we propose a novel approach to construct compact representations of model abilities as shown in Figure 1, drawing inspiration from the item response theory (IRT). IRT is a well-established statistical framework used in education and psychology to measure latent abilities through standardized tests. It models the interactions between individuals’ performance, their latent abilities, and query characteristics, such as difficulty and discrimination (as in the 2-parameter IRT). By treating queries as tests and LLMs as test respondents, we analogously model the likelihood of an LLM  $m$  correctly answering a query  $q$  as a mathematical function of  $m$ ’s latent abilities  $\theta_m$ , query

<sup>1</sup>For example, NeurIPS 2025’s dataset and benchmark track received more than 1,900 submissions, and in 2024, more than 1,200 submissions; a large proportion of accepted submissions are LLM benchmarks.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

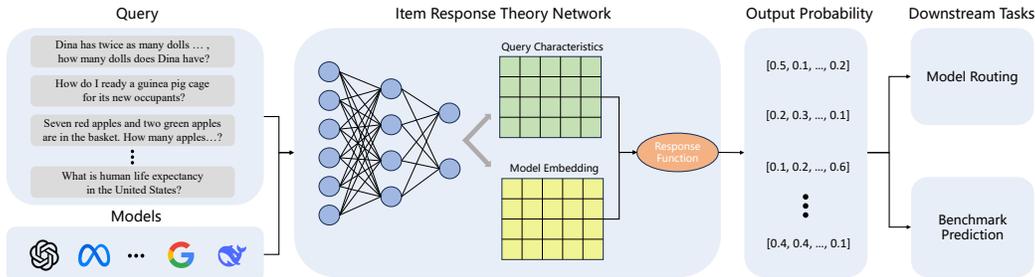


Figure 1: Overview of the IrtNet framework for learning LLM representations. IrtNet learns model embeddings based on models’ past query answering performance and outputs probabilities that models answer correctly. The output probability can be directly applied to downstream tasks containing model routing and benchmark prediction.

$q$ ’s difficulty  $\beta_q$ , and  $q$ ’s discrimination power  $\alpha_q$  between models. To estimate the model’s compact representation of abilities ( $\theta_m$ ) and the query characteristics ( $\beta_q$  and  $\alpha_q$ ), we present a neural network *IrtNet*. This architecture enables end-to-end training, allowing for the simultaneous optimization of model latent abilities and query characteristics to align with the ground truth of whether a given model correctly answers a given query.

Comprehensive experiments show that our approach not only achieves state-of-the-art performance on downstream tasks but also produces interpretable learned representations. Specifically, in model routing, our method achieves an average accuracy of 67.4%, significantly outperforming strong baselines like RouterDC (Chen et al., 2024) (54.9%), MODEL-SAT (Zhang et al., 2025c) (56.7%), Avengers-Pro (Zhang et al., 2025e) (62.1%), EmbedLLM (Zhuang et al., 2025) (60.2%). In benchmark prediction, our method demonstrates remarkable data efficiency. It reaches 69.9% accuracy using less than 4% of the training data, matching the state-of-the-art (EmbedLLM) performance achieved with the full training set. Moreover, our analysis validates that the learned parameters encode meaningful information: models’ compact representations form clear clusters in the latent space based on model family (e.g., Llama, Qwen series) and specialization (e.g., models trained for coding and mathematics); the learned query difficulty parameter  $\beta_q$  exhibits a near-perfect negative correlation with true benchmark scores, with a Pearson correlation coefficient of -0.9721. These findings indicate that IrtNet successfully captures the intrinsic model abilities and query characteristics, providing a powerful new tool for evaluation, selection, and management of the vast, rapidly expanding LLM ecosystem.

## 2 PRELIMINARY

We denote the compact representation of an LLM  $m$  by a  $d$ -dimensional embedding vector  $\theta_m \in \mathbb{R}^d$ , which encodes the model’s strengths and weaknesses. Let  $M = \{m_1, m_2, \dots, m_n\}$  be a set of  $n$  distinct LLMs, and let  $Q = \{q_1, q_2, \dots, q_k\}$  be a set of  $k$  distinct queries. For any given model-query pair  $(m, q)$ , where  $m \in M$  and  $q \in Q$ , the model’s answer to the query can be represented by a binary outcome  $y \in \{0, 1\}$ , where  $y = 1$  signifies a correct answer and  $y = 0$  signifies an incorrect one.

Our objective is to learn a probability mass function  $f_\theta(m, q) = \Pr(y = 1|m, q)$ , parameterized by  $\theta = (\theta_1, \dots, \theta_n)$ , such that for any given model-query pair, the function  $f_\theta$  can accurately predict whether the model can answer the query correctly or not. Once learned, the function  $f_\theta$  can facilitate several important downstream tasks:

**Model Routing.** Given a set of candidate models, model routing aims to assign each query to the most suitable model, avoiding the need for every model to answer it exhaustively Shnitzer et al. (2023). In this way, the ensemble of models effectively leverages collective intelligence, allowing the group to solve a broader range of tasks that any individual model alone cannot accomplish. Formally, let  $\mathcal{M} \subseteq M$  be the set of candidate models; based on the function  $f_\theta$ , for any query  $q$ , model routing can be achieved by selecting the model  $m^*$  that yields the highest probability of

generating a correct answer, i.e.,

$$m^* = \arg \max_{m \in \mathcal{M}} f_\theta(m, q). \quad (1)$$

**Benchmark Prediction.** As exhaustive LLM evaluation is compute-intensive and time-consuming, benchmark prediction seeks to estimate the overall performance of an LLM from only a small subset of evaluation data (Vivek et al., 2023). Formally, let  $\Omega$  be the set of queries that are not included when learning the function  $f_\theta$ , i.e.,  $\Omega \cap Q = \emptyset$ , where typically  $|\Omega| \gg |Q|$ ; for a given model  $m$ , benchmark prediction generates a predicted accuracy  $\hat{S}$  for the set  $\Omega$  by feeding each query from this set into the function  $f_\theta$ , which outputs the predicted probability of generating correct answers to these queries, i.e.,

$$\hat{S} = \frac{1}{|\Omega|} \sum_{q \in \Omega} f_\theta(m, q). \quad (2)$$

In practice, as  $f_\theta$  is a neural predictor, a single pass through all queries in the set  $\Omega$  can yield the predicted accuracy for all models in the set  $M$ . This avoids the need to run multiple passes (one pass for each model) as in traditional benchmarking LLMs, thereby allowing for efficient, scalable LLM evaluation.

Overall, establishing compact representations of LLMs provides a unified framework for understanding a model’s strengths and weaknesses, supporting a range of downstream applications that have attracted significant recent interest.

### 3 METHODOLOGY

#### 3.1 MODELING LLM ABILITIES VIA ITEM RESPONSE THEORY

The item response theory (IRT) is a well-established statistical framework to measure the latent abilities of respondents through testings (Cai et al., 2016). It is widely applied in education and psychology, and underpins the scoring scales of high-stakes exams such as the GRE and GMAT.

The IRT is based on the idea that the probability of a correct response to an item (or a question) is a mathematical function of a person’s latent traits, indicating the person’s ability, as well as the item parameters. Here, we focus on the two-parameter IRT model, which considers an item’s difficulty and discrimination (its ability to differentiate between individuals) (Reise & Waller, 2009).

Understanding LLM abilities through their performance on queries is analogous to assessing a person’s abilities based on their performance on standardized tests. Thus, inspired by the IRT, we treat each LLM  $m$  as a respondent with a latent trait, and each query  $q$  as an item with two parameters that characterize its difficulty and discriminative power.

Formally, let  $\alpha_q \in \mathbb{R}^d$  denote the query  $q$ ’s  $d$ -dimensional discrimination parameter, and  $\beta_q \in \mathbb{R}$  denote  $q$ ’s difficulty parameter; we assume that the function  $f_\theta(m, q)$ , which predicts whether the model  $m$  can correctly answer the query  $q$ , takes the following form of an item response function:

$$f_\theta(m, q) = \sigma(\alpha_q^\top \theta_m - \beta_q) = \frac{1}{1 + e^{-(\alpha_q^\top \theta_m - \beta_q)}}, \quad (3)$$

where  $\theta_m \in \mathbb{R}^d$  is the model  $m$ ’s compact representation (defined in Section 2), indicating its latent abilities, and  $\sigma(\cdot)$  is the logistic link function. The discrimination parameter  $\alpha_q$  suggests how important each latent ability dimension for the query is, where a higher value indicates that proficiency in that dimension is more critical for a correct answer. The dot product  $\alpha_q^\top \theta_m$  computes how well the model’s abilities align with the discriminative power of the query. Essentially, it measures the model’s fit to the query in terms of both the model’s ability in relevant latent dimensions (encoded by  $\theta_m$ ) and the importance of those dimensions to the query (encoded by  $\alpha_q$ ).

#### 3.2 IRTNET: LEARNING COMPACT REPRESENTATIONS

We propose to jointly learn models’ compact representations  $\theta = (\theta_1, \dots, \theta_n)$ , queries’ discrimination parameters  $\alpha = (\alpha_1, \dots, \alpha_k)$ , and difficulty parameters  $\beta = (\beta_1, \dots, \beta_k)$  through a neural network *IrtNet* in an end-to-end manner.

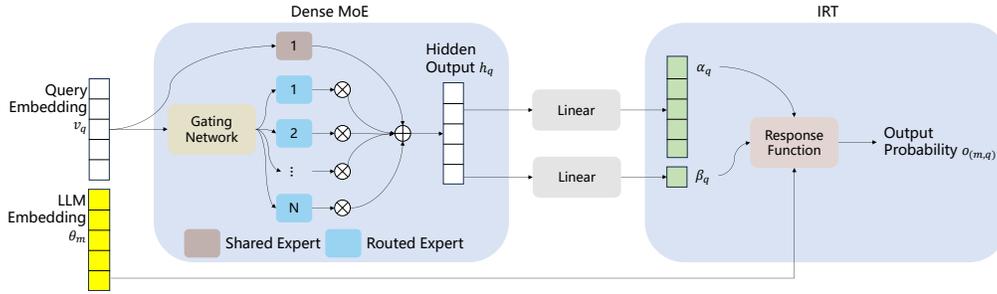


Figure 2: The architecture of IrtNet. A query embedding is processed through a dense MoE layer and subsequent linear layers to generate the query’s discrimination  $\alpha_q$  and difficulty  $\beta_q$  parameters. These parameters are then combined with an LLM embedding  $\theta_m$  via the response function to compute the final output probability.

We illustrate the architecture of IrtNet in Figure 2. First, each query  $q$  is converted into a semantic embedding  $v_q$  using a pre-trained embedding model. Next,  $v_q$  is fed into the MoE layer of the IrtNet. We use a dense MoE layer which always activates all  $N$  experts, with an auxiliary-loss-free load balancing strategy (Liu et al., 2024) to promote a diverse weight distribution. This design aims to capture diverse, multi-faceted understanding of the query, improving prediction accuracy while avoiding the training instability problem often encountered in sparse MoE layers. The query’s hidden output  $h_q$  is obtained by

$$h_q = \text{MoE}(v_q) = \text{SharedExpert}(v_q) + \sum_{i=1}^N w_i \cdot \text{RoutedExpert}_i(v_q). \quad (4)$$

Next, we employ two independent linear layers to obtain the discrimination parameter  $\alpha$  and difficulty parameter  $\beta$ :

$$\alpha_q = \text{Linear}(h_q, d), \quad \beta_q = \text{Linear}(h_q, 1). \quad (5)$$

where  $d$  is the dimension of  $\theta_m$ . Finally, the IrtNet combines the model embedding  $\theta_m$  with the query characteristic parameters  $\alpha_q$  and  $\beta_q$  to compute the output  $o_{(m,q)}$  based on the response function shown in Equation 3.

To learn these parameters, we define the objective as minimizing the discrepancy between the predicted probabilities  $o_{(m,q)}$  and the ground-truth labels  $y$  across the entire training dataset. Specifically, for a training set  $\mathcal{D}$  consisting of samples  $(m, q, y)$ , the overall loss  $\mathcal{L}$  is formulated as the sum of the binary cross-entropy losses for all samples:

$$\mathcal{L} = - \sum_{(m,q,y) \in \mathcal{D}} [y \log(o_{(m,q)}) + (1 - y) \log(1 - o_{(m,q)})]. \quad (6)$$

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

**Datasets** We use the same data as EmbedLLM (Zhuang et al., 2025). Data repository: <https://huggingface.co/datasets/RZ412/EmbedLLM>. We applied a majority vote to consolidate multiple answers from a model to the same query. This step ensures a unique ground truth for each model-query pair, which is especially critical for the test set. The datasets contain 35,673 queries from 10 public benchmarks, including ASDiv (Miao et al., 2020), GPQA (Rein et al., 2024), GSM8K (Cobbe et al., 2021), MathQA (Amini et al., 2019), LogiQA (Liu et al., 2020), MedMCQA (Pal et al., 2022), MMLU (Hendrycks et al., 2021), SocialIQA (Sap et al., 2019), PIQA (Bisk et al., 2020), and TruthfulQA (Lin et al., 2022). The correctness of answers from 112 open-source language models to those queries was evaluated. The queries were converted into 768-dimensional embeddings using the all-mpnet-base-v2 (Reimers & Gurevych, 2019) sentence transformer. The queries were split into a training set of 29,673 queries, a validation set of 3,000 queries, and a test set of 3,000 queries. LLMs contained are shown in Appendix A.6

**Hyperparameters** The dimension  $d$  of the compact LLM representation  $\theta_m$  is 232. The number of experts in the MoE layer is set to 40, which is 4 times the number of datasets.

## 4.2 MODEL ROUTING

In this section, we test IrtNet’s ability to serve as an effective and intelligent router in a multi-model environment. The objective is to leverage the fine-grained predictions from our framework to select the best model from a diverse pool to handle a given query, thereby maximizing both accuracy and efficiency. For any given query  $q$ , we compute the predicted success probability  $P(y = 1|m, q)$  for all candidate models and route the query to the model with the highest probability.

Table 1: Model routing accuracy (%) comparison. **Micro** denotes the micro-averaged accuracy, calculated on the total correct predictions across all datasets. **Macro** denotes the macro-averaged accuracy, computed by first calculating the accuracy for each benchmark individually and then averaging these benchmark scores. We use **bold** to indicate the best results.

Method	ASD	GPQA	GSM	Math	Logi	Med	Mmlu	Soci	PIQA	Tru	Micro	Macro
RouterDC	62.1	21.2	70.5	40.5	41.2	54.5	70.9	30.9	80.6	50.0	54.9	52.2
EmbedLLM	34.9	24.8	82.1	47.7	31.4	65.8	80.1	<b>37.0</b>	86.6	<b>52.7</b>	60.2	54.3
MODEL-SAT	6.57	22.6	80.4	43.0	47.0	62.4	80.8	29.6	83.6	35.1	56.7	49.1
Avengers-Pro	12.1	<b>29.4</b>	<b>89.3</b>	55.3	<b>49.0</b>	72.9	83.2	30.9	86.6	41.9	62.1	55.7
IrtNet (Ours)	<b>66.7</b>	28.6	<b>89.3</b>	<b>57.8</b>	<b>49.0</b>	<b>74.0</b>	<b>86.2</b>	34.0	<b>87.3</b>	47.3	<b>67.4</b>	<b>62.0</b>

We compare our method with four advanced routing methods, RouterDC (Chen et al., 2024), EmbedLLM (Zhuang et al., 2025), MODEL-SAT (Zhang et al., 2025c) (with Qwen3-0.6B-Base (Yang et al., 2025) as base model), and Avengers-Pro (Zhang et al., 2025d). As shown in Table 1, our method achieves the best performance on both micro-average and macro-average metrics across 10 benchmarks. Specifically, IrtNet achieves a final micro-average accuracy of 67.4%, significantly outperforming strong baselines like EmbedLLM (60.2%) and Avengers-Pro (62.1%), which showcases its immense potential for the model routing task.

## 4.3 BENCHMARK PREDICTION

We use two settings to test the benchmark prediction ability of IrtNet: in-distribution (ID) and out-of-distribution (OOD). In the ID setting, while the two sets  $\Omega$  and  $Q$  have no overlap, they may still be closely related. In the OOD case, performance on a dataset may be predicted based on the model’s performance on a different dataset.

**ID Correctness Prediction.** In this setting, we evaluate IrtNet’s core capability for correctness prediction on model-query pairs. The results, presented in Table 2, highlight our model’s exceptional data efficiency and predictive power.

Table 2: Correctness prediction accuracy (%) on different training data sizes.

method	Dataset Size						
	1K	5K	10K	15K	20K	25K	Full (29K)
KNN	62.6	63.0	64.4	65.1	64.6	64.4	64.7
EmbedLLM	60.8	64.2	66.5	67.9	69.1	69.9	70.6
IrtNet (ours)	<b>69.9</b>	<b>71.5</b>	<b>71.7</b>	<b>71.8</b>	<b>72.0</b>	<b>72.1</b>	<b>72.2</b>

A key advantage of our framework is its ability to achieve strong performance with remarkably little training data. With just 1,000 queries—less than 4% of the full training set—IrtNet achieves a prediction accuracy of 69.9%. This performance significantly surpasses both the traditional KNN baseline (62.6%) and the state-of-the-art EmbedLLM (60.8%) trained on the same amount of data. Notably, our model’s accuracy with only 1K queries already approaches the performance of EmbedLLM trained on the entire dataset (70.6%).

As the training set size increases, IrtNet consistently outperforms other baselines, reaching a final accuracy of 72.2% on the full training set. This demonstrates that our framework not only learns rapidly from limited samples but also scales effectively. The results strongly suggest that by modeling the interaction between model abilities and query characteristics, IrtNet captures their complex relationship more efficiently and accurately than existing methods.

**OOD Benchmark Prediction.** In this setting, we test IrtNet’s ability to generalize its learned representations to make macroscopic performance predictions. The objective is to forecast an LLM’s accuracy score on an entire benchmark it has never seen during training. This experiment uses a leave-one-out approach: we train IrtNet on all data except for one target benchmark and then use the trained model to predict the overall accuracy of all LLMs on that held-out benchmark. We exclude ASDiv and SocialIQA from this analysis, as the near-uniform scores across all models (approximately 0 and 0.3, respectively) suggest the data might represent noise rather than meaningful performance variation. This OOD setup serves as the ultimate test of whether our model representations can truly understand and generalize the abstract concept of LLM abilities.

Table 3: Root mean square error (RMSE) for OOD benchmark prediction. **Overall** denotes the total prediction error, calculated by treating all benchmarks as a single, unified test set.

Method	GPQA	GSM	Math	Logi	Med	Mmlu	PIQA	Tru	Overall
EmbedLLM	0.26	0.20	0.09	0.11	0.04	0.25	0.40	<b>0.11</b>	0.21
IrtNet (Ours)	0.26	<b>0.19</b>	0.09	0.11	0.04	<b>0.19</b>	<b>0.36</b>	0.12	<b>0.19</b>

The results in Table 3 demonstrate IrtNet’s robust benchmark accuracy prediction. Across the eight benchmarks, IrtNet achieves a lower root mean square error (RMSE) than EmbedLLM on three datasets and performs equally on four. This leads to an overall RMSE of 0.19, representing a nearly 10% error reduction over EmbedLLM’s 0.21.

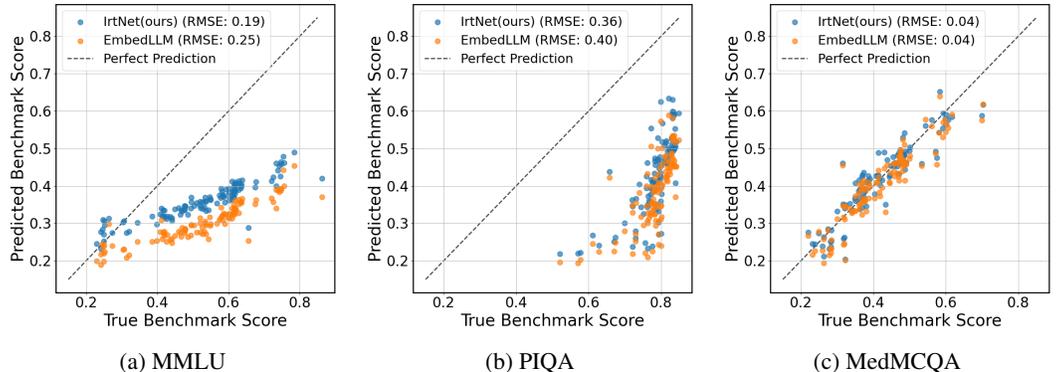


Figure 3: Predicted vs. true benchmark scores (in [0-1]) on three OOD benchmarks. The scatter plots represent the predicted LLM scores by IrtNet and EmbedLLM. IrtNet’s predictions (blue dots) align more closely with the perfect prediction diagonal line on MMLU and PIQA, which means lower prediction errors. IrtNet and EmbedLLM are tied on MedMCQA with the predicted scores almost coinciding with the true scores.

Figure 3 presents the specific prediction distributions on three representative benchmarks (MMLU, PIQA, and MedMCQA), vividly illustrating how IrtNet’s predictions achieve a smaller prediction error compared to EmbedLLM or cluster very tightly around the true scores.

#### 4.4 ABLATION STUDY (REVISED)

In this section, we investigate the necessity of both the MoE architecture and the IRT formulation. We ablate the MoE layer by replacing it with an MLP network of the same number of parameters and ablate the IRT formulation by replacing it with logistic regression.

As shown in Table 4, after removing the MoE layer or IRT formulation, IrtNet exhibits a significant performance drop in the model routing task and a slight decline in the correctness prediction task.

Table 4: Ablation study on the MoE layer and IRT formulation of IrtNet. **Routing** denotes the model routing task and **Correctness** denotes the correctness prediction task.

Task	IrtNet	w/o MoE	w/o IRT	w/o Both (EmbedLLM)
<b>Routing</b>	67.4	64.0 $\downarrow$ 3.4	63.8 $\downarrow$ 3.6	60.2 $\downarrow$ 6.2
<b>Correctness</b>	72.2	71.3 $\downarrow$ 0.9	71.1 $\downarrow$ 1.1	70.6 $\downarrow$ 1.6

This ablation result fully demonstrates the effectiveness of the MoE structure and the IRT formulation. Furthermore, the more substantial decline in the model routing task indicates that IrtNet is particularly effective at understanding the relative ranking among models. This is likely because IrtNet produces a more discriminative  $\alpha_q$  representation, which more accurately captures the distinctions among model abilities.

## 5 INTERPRETABILITY ANALYSIS

### 5.1 UNDERSTANDING DISCRIMINATION

In our framework, the discrimination vector  $\alpha_q$  provides a rich, quantitative profile of a query’s intrinsic characteristic. A high value in a particular dimension of  $\alpha_q$  signifies that the query places a strong demand on the corresponding latent ability, making a model’s proficiency in that dimension more critical for a correct answer. Thus, the discrimination vector reflects which abilities a query is designed to test and how strongly it tests them. To validate this interpretation, we conduct a visualization experiment on the test set. We use the trained IrtNet to compute the discrimination vectors for queries on the test set and then project these high-dimensional vectors into a two-dimensional space using t-SNE.

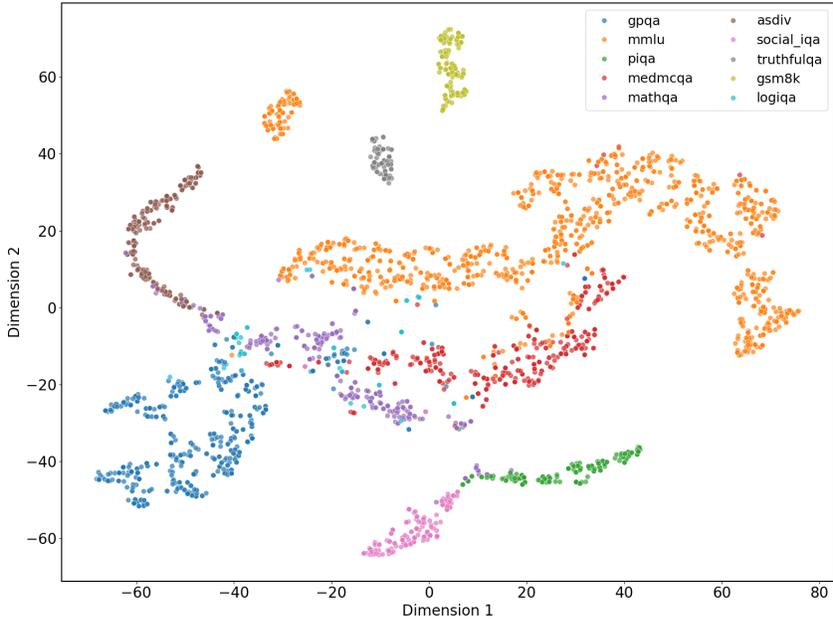


Figure 4: T-SNE visualization of learned query discrimination vectors  $\alpha_q$ .

Remarkably, despite IrtNet never being exposed to any dataset labels during training, the resulting visualization in Figure 4 reveals a clear and well-defined clustering structure. We observe that queries originating from the same benchmark naturally form distinct semantic groups in this learned space. This provides strong evidence that the learned discrimination vector  $\alpha_q$  has successfully captured the unique demands of different query types. The spontaneous emergence of these clusters demonstrates that IrtNet has effectively modeled the distinct discriminative properties inherent to each semantic query group, mapping them into an interpretable space.

### 5.2 VALIDATING DIFFICULTY

We perform a quantitative analysis to validate that our learned difficulty parameter  $\beta_q$  is a meaningful measure of a query’s intrinsic challenge. For our experimental setup, we first establish an objective difficulty standard for each benchmark. This standard is defined as the average ground-truth accuracy achieved by the entire pool of 112 models on that dataset. A lower average score naturally corresponds to a higher objective difficulty. We then compare this objective standard against the average learned difficulty parameter  $\beta_q$  for each respective benchmark.

Table 5: Comparison of average model accuracy (in [0, 1]) and learned difficulty  $\beta_q$  across benchmarks.

	ASD	GPQA	Logi	Math	Soci	Tru	Med	GSM	Mmlu	PIQA
Average Accuracy	0.04	0.21	0.29	0.33	0.34	0.36	0.42	0.42	0.53	0.78
Average $\beta_q$	1.69	0.71	0.45	0.40	0.33	0.32	0.16	0.17	-0.08	-0.78

As shown in Table 5, our findings reveal an exceptionally strong negative correlation between the two: as the average model accuracy on a benchmark decreases (i.e., it gets harder), the learned  $\beta_q$  value consistently increases. Quantitatively, the Pearson correlation coefficient between the objective average accuracy and the learned average difficulty  $\beta_q$  is -0.9721. This near-perfect correlation provides compelling evidence that the difficulty parameter learned by IrtNet is not an arbitrary value, but a valid and reliable metric that accurately captures the empirical hardness of the queries.

### 5.3 PROBING LLM EMBEDDING

To investigate whether the learned LLM embedding  $\theta_m$  captures meaningful model characteristics, we conduct a similarity validation experiment. Our hypothesis is that if the vectors are coherent, models sharing fundamental traits should be geometrically closer in the learned space. To test this, we partition our pool of models into distinct groups based on two criteria: model family, such as the Qwen and Llama families, and domain specialization for tasks like medicine, code, and math. We then calculate and compare two metrics: the average intra-community L2 distance, which measures the distance between models within a single group, and the average inter-community L2 distance, which measures the distance from that group’s models to all outside models.

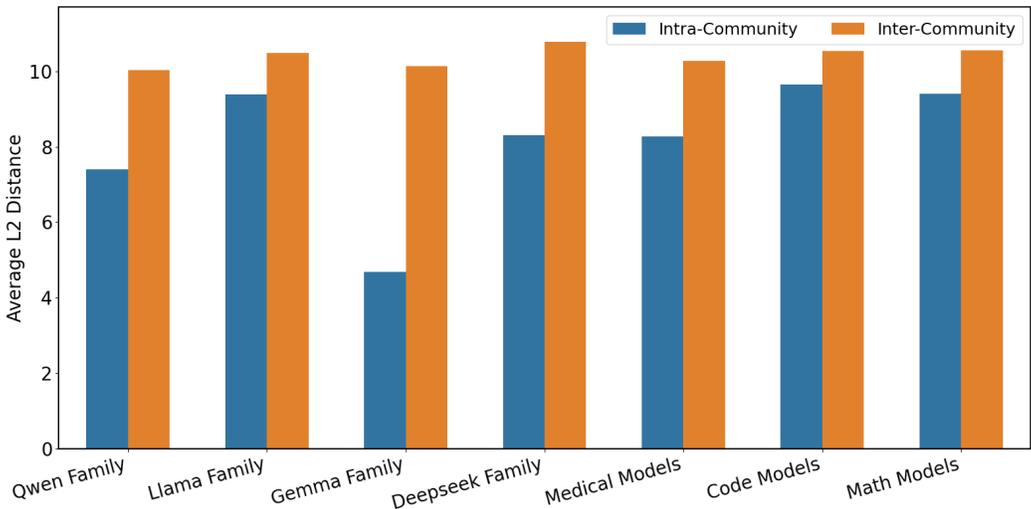


Figure 5: Comparison of intra-community and inter-community L2 distances for LLM embeddings.

Figure 5 provides compelling evidence for our hypothesis. Across all defined groups, the average intra-community distance is consistently and significantly smaller than the average inter-community distance. This clear geometric clustering holds true for models grouped by both architectural lineage, such as the Llama and Gemma families, or by functional specialization, like the groups of Code and

432 Math models. This finding demonstrates that model embedding  $\theta_m$  is a meaningful, well-structured  
433 representation that effectively encodes a model’s specialized abilities.  
434

## 435 6 RELATED WORK 436

437 **Representation Learning** The concept of learning compact vector representations for complex,  
438 high-dimensional objects is a well-established paradigm in machine learning . For instance, mod-  
439 els like SentenceTransformer (Reimers & Gurevych, 2019) and Qwen3-Embedding (Zhang et al.,  
440 2025b) effectively embed sentences into a dense embedding space; in knowledge graphs, methods  
441 like TransE (Bordes et al., 2013) learn low-dimensional embeddings for entities and relations. More  
442 recently, EmbedLLM (Zhuang et al., 2025) generalizes this paradigm to model LLM abilities with  
443 an encoder-decoder approach, demonstrating that learning a compact representation of LLM abil-  
444 ities can facilitate multiple downstream tasks. Our work is most closely related to EmbedLLM,  
445 but bears key conceptual and methodological differences. We apply IRT to model LLM abilities,  
446 providing a theory-driven approach to this paradigm, and introduce an MoE-based method to ex-  
447 plicitly learn query characteristics (discrimination and difficulty) that EmbedLLM does not capture.  
448 Furthermore, our method leads to significant outperformance compared to EmbedLLM.  
449

450 **Model Routing** Model routing for LLMs (Shnitzer et al., 2023; Lu et al., 2024; Srivatsa et al.,  
451 2024) has emerged as a critical strategy for efficiently managing a diverse suite of models. Recent  
452 research in this field generally falls into two streams. The first emphasizes striking a balance be-  
453 tween performance and computational efficiency (Jiang et al., 2023; Ong et al., 2024; Feng et al.,  
454 2024; Wang et al., 2025; Zhang et al., 2025d). The second prioritizes pushing model performance  
455 to the highest possible levels (Lu et al., 2024; Chen et al., 2024; Zhang et al., 2025c;e). Roun-  
456 terDC (Chen et al., 2024) proposes a dual-contrastive learning framework to better align queries and  
457 model representations. Model-SAT (Zhang et al., 2025c) creates a capability representation for each  
458 model through a lightweight aptitude test to select the most suitable model. The Avengers (Zhang  
459 et al., 2025e) presents a training-free, clustering-based routing framework that selects the optimal  
460 model. While our method does not specifically target model routing, the router developed using our  
461 approach significantly outperforms state-of-the-art routing methods.

462 **Benchmark Prediction** Evaluating LLMs on large benchmarks requires substantial computa-  
463 tional and financial cost, which has motivated research on benchmark prediction that seeks to es-  
464 timate overall performance from only a subset of representative data. Core-set selection methods  
465 such as Anchor Points (Vivek et al., 2023) attempt to identify a small number of informative queries  
466 that preserve model rankings nearly as well as full benchmarks. Other work has shown that even  
467 simple random subsets combined with regression models can provide surprisingly strong estimates  
468 of average accuracy (Zhang et al., 2025a). More recently, IRT-based methods such as tinyBench-  
469 marks (Polo et al., 2024b) have demonstrated that psychometric modeling can reduce evaluation  
470 costs even further. While our approach to learning representations of model abilities can be applied  
471 to benchmark prediction, it differs from the above studies that focus on data selection strategies. Dis-  
472 tinct from these, observational scaling laws (Polo et al., 2024a; Ruan et al., 2024) predict aggregate  
473 benchmark performance by modeling latent skills from model metadata. Unlike these macro-level  
474 approaches, we leverage query-level semantic prediction and then aggregate the results across the  
475 entire benchmark.

## 476 7 CONCLUSION 477

478 In this paper, we pioneer the application of item response theory (IRT) to formally model LLM  
479 abilities. We introduce IrtNet, a IRT-based framework that learns compact representations of LLM  
480 abilities. Based on a Mixture-of-Experts architecture, it jointly learns model embeddings alongside  
481 query difficulty and discrimination. Extensive Experiments demonstrate that IrtNet sets a new state-  
482 of-the-art in model routing and achieves highly data-efficient, more accurate benchmark prediction.  
483 Furthermore, the learned representations are also interpretable. The difficulty parameter strongly  
484 correlates with empirical results, and that model embeddings naturally cluster by family and func-  
485 tion. Overall, IrtNet provides a robust and insightful tool for effective model evaluation, selection,  
and analysis in the growing LLM ecosystem.

486 REPRODUCIBILITY STATEMENT  
487

488 To facilitate the reproduction of our results, we have provided the implementation of our method in  
489 the supplementary materials. The datasets used in this study are all publicly accessible, and details  
490 can be found in the experiments section. Furthermore, to contribute to the open-source community,  
491 we pledge to release a cleaned, user-friendly, and well-documented version of our source code on a  
492 public repository (e.g., GitHub) upon acceptance of the manuscript.  
493

494 REFERENCES  
495

- 496 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh  
497 Hajishirzi. Mathqa: Towards interpretable math word problem solving with operation-based for-  
498 malisms. In *Proceedings of the 2019 Conference of the North American Chapter of the Associ-  
499 ation for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Min-  
500 neapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, 2019.
- 501 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
502 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
503 models. *arXiv preprint arXiv:2108.07732*, 2021.  
504
- 505 Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. PIQA: reasoning about  
506 physical commonsense in natural language. In *The Thirty-Fourth AAAI Conference on Artificial  
507 Intelligence, AAAI*, 2020.
- 508 Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko.  
509 Translating embeddings for modeling multi-relational data. *Advances in neural information pro-  
510 cessing systems*, 26, 2013.  
511
- 512 Mats Byrkjeland, Frederik Gørvell de Lichtenberg, and Björn Gambäck. Ternary twitter sentiment  
513 classification with distant supervision and sentiment-specific word embeddings. In *Proceedings  
514 of the 9th workshop on computational approaches to subjectivity, sentiment and social media  
515 analysis*, pp. 97–106, 2018.
- 516 Li Cai, Kilchan Choi, Mark Hansen, and Lauren Harrell. Item response theory. volume 3, pp.  
517 297–321. *Annual Reviews*, 2016.  
518
- 519 Shuhao Chen, Weisen Jiang, Baijiong Lin, James T. Kwok, and Yu Zhang. Routerdc: Query-based  
520 router by dual contrastive learning for assembling large language models. In *Advances in Neu-  
521 ral Information Processing Systems 38: Annual Conference on Neural Information Processing  
522 Systems, NeurIPS*, 2024.  
523
- 524 Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema  
525 Moussa, Matt Beane, Ting-Hao Huang, Bryan R Routledge, et al. Finqa: A dataset of numerical  
526 reasoning over financial data. In *Proceedings of the 2021 Conference on Empirical Methods in  
527 Natural Language Processing*, pp. 3697–3711, 2021.
- 528 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
529 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
530 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.  
531
- 532 Mucong Ding, Chenghao Deng, Jocelyn Choo, Zichu Wu, Aakriti Agrawal, Avi Schwarzschild,  
533 Tianyi Zhou, Tom Goldstein, John Langford, Animashree Anandkumar, et al. Easy2hard-bench:  
534 Standardized difficulty labels for profiling llm performance and generalization. *Advances in Neu-  
535 ral Information Processing Systems*, 37:44323–44365, 2024.
- 536 Tao Feng, Yanzhen Shen, and Jiaxuan You. Graphrouter: A graph-based router for llm selections.  
537 *arXiv preprint arXiv:2410.03834*, 2024.  
538
- 539 Evan Frick, Connor Chen, Joseph Tennyson, Tianle Li, Wei-Lin Chiang, Anastasios N Angelopou-  
los, and Ion Stoica. Prompt-to-leaderboard. *arXiv preprint arXiv:2502.14855*, 2025.

- 540 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
541 Steinhardt. Measuring massive multitask language understanding. In *9th International Confer-*  
542 *ence on Learning Representations, ICLR*, 2021.
- 543 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando  
544 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free  
545 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- 546 Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. Llm-blender: Ensembling large language models  
547 with pairwise ranking and generative fusion. In *Proceedings of the 61st Annual Meeting of the*  
548 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14165–14178, 2023.
- 549 Wittawat Jitkrittum, Harikrishna Narasimhan, Ankit Singh Rawat, Jeevesh Juneja, Congchao Wang,  
550 Zifeng Wang, Alec Go, Chen-Yu Lee, Pradeep Shenoy, Rina Panigrahy, et al. Universal model  
551 routing for efficient llm inference. *arXiv preprint arXiv:2502.08773*, 2025.
- 552 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan  
553 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth*  
554 *International Conference on Learning Representations*, 2023.
- 555 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human  
556 falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational*  
557 *Linguistics, ACL*, 2022.
- 558 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,  
559 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*  
560 *arXiv:2412.19437*, 2024.
- 561 Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A  
562 challenge dataset for machine reading comprehension with logical reasoning. In *Proceedings of*  
563 *the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI*, 2020.
- 564 Keming Lu, Hongyi Yuan, Runji Lin, Junyang Lin, Zheng Yuan, Chang Zhou, and Jingren Zhou.  
565 Routing to the expert: Efficient reward-guided ensemble of large language models. In *Proceedings*  
566 *of the 2024 Conference of the North American Chapter of the Association for Computational*  
567 *Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico*  
568 *City, Mexico, June 16-21, 2024*, pp. 1964–1974. Association for Computational Linguistics, 2024.
- 569 Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing  
570 english math word problem solvers. In *Proceedings of the 58th Annual Meeting of the Association*  
571 *for Computational Linguistics, ACL*, 2020.
- 572 Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez,  
573 M Waleed Kadous, and Ion Stoica. Routellm: Learning to route llms from preference data.  
574 In *The Thirteenth International Conference on Learning Representations*, 2024.
- 575 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale  
576 multi-subject multi-choice dataset for medical domain question answering. In *Conference on*  
577 *Health, Inference, and Learning, CHIL*, 2022.
- 578 Zhihong Pan, Kai Zhang, Yuze Zhao, and Yupeng Han. Route to reason: Adaptive routing for llm  
579 and reasoning strategy selection. *arXiv preprint arXiv:2505.19435*, 2025.
- 580 Felipe Maia Polo, Seamus Somerstep, Leshem Choshen, Yuekai Sun, and Mikhail Yurochkin. Sloth:  
581 scaling laws for llm skills to predict multi-benchmark performance across families. *arXiv preprint*  
582 *arXiv:2412.06540*, 2024a.
- 583 Felipe Maia Polo, Lucas Weber, Leshem Choshen, Yuekai Sun, Gongjun Xu, and Mikhail  
584 Yurochkin. tinybenchmarks: evaluating llms with fewer examples. 2024b.
- 585 Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada  
586 Mihalcea. Meld: A multimodal multi-party dataset for emotion recognition in conversations.  
587 In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pp.  
588 527–536, 2019.

- 594 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-  
595 networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*  
596 *Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-*  
597 *IJCNLP*, 2019.
- 598  
599 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di-  
600 rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench-  
601 mark. In *First Conference on Language Modeling*, 2024.
- 602  
603 Steven P Reise and Niels G Waller. Item response theory and clinical measurement. volume 5, pp.  
604 27–48. *Annual Reviews*, 2009.
- 605  
606 Yangjun Ruan, Chris J Maddison, and Tatsunori B Hashimoto. Observational scaling laws and  
607 the predictability of language model performance. *Advances in Neural Information Processing*  
608 *Systems*, 37:15841–15892, 2024.
- 609  
610 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Socialiqa: Common-  
611 sense reasoning about social interactions. *CoRR*, abs/1904.09728, 2019.
- 612  
613 Tal Shnitzer, Anthony Ou, Mírian Silva, Kate Soule, Yuekai Sun, Justin Solomon, Neil Thompson,  
614 and Mikhail Yurochkin. Large language model routing with benchmark datasets. *arXiv preprint*  
615 *arXiv:2309.15789*, 2023.
- 616  
617 KV Srivatsa, Kaushal Kumar Maurya, and Ekaterina Kochmar. Harnessing the power of multiple  
618 minds: Lessons learned from llm routing. *arXiv preprint arXiv:2405.00467*, 2024.
- 619  
620 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
621 Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. Challenging big-bench tasks and  
622 whether chain-of-thought can solve them. In *Findings of the Association for Computational Lin-*  
623 *guistics: ACL 2023*, pp. 13003–13051, 2023.
- 624  
625 Rajan Vivek, Kawin Ethayarajh, Diyi Yang, and Douwe Kiela. Anchor points: Benchmarking  
626 models with much fewer examples. 2023.
- 627  
628 Xinyuan Wang, Yanchi Liu, Wei Cheng, Xujiang Zhao, Zhengzhang Chen, Wenchao Yu, Yanjie Fu,  
629 and Haifeng Chen. Mixllm: Dynamic routing in mixed large language models. In *Proceedings of*  
630 *the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational*  
631 *Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 10912–10922, 2025.
- 632  
633 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming  
634 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-  
635 task language understanding benchmark. *Advances in Neural Information Processing Systems*,  
636 37:95266–95290, 2024.
- 637  
638 Chulin Xie, Yangsibo Huang, Chiyuan Zhang, Da Yu, Xinyun Chen, Bill Yuchen Lin, Bo Li, Badih  
639 Ghazi, and Ravi Kumar. On memorization of large language models in logical reasoning. *arXiv*  
640 *preprint arXiv:2410.23123*, 2024.
- 641  
642 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
643 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*  
644 *arXiv:2505.09388*, 2025.
- 645  
646 Guanhua Zhang, Florian E Dorner, and Moritz Hardt. How benchmark prediction from fewer data  
647 misses the mark. 2025a.
- 648  
649 Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie,  
650 An Yang, Dayiheng Liu, Junyang Lin, et al. Qwen3 embedding: Advancing text embedding and  
651 reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025b.
- 652  
653 Yi-Kai Zhang, De-Chuan Zhan, and Han-Jia Ye. Capability instruction tuning: A new paradigm for  
654 dynamic llm routing. *arXiv preprint arXiv:2502.17282*, 2025c.

648 Yiqun Zhang, Hao Li, Jianhao Chen, Hangfan Zhang, Peng Ye, Lei Bai, and Shuyue Hu. Be-  
649 yond gpt-5: Making llms cheaper and better via performance-efficiency optimized routing. *arXiv*  
650 *preprint arXiv:2508.12631*, 2025d.

651 Yiqun Zhang, Hao Li, Chenxu Wang, Linyao Chen, Qiaosheng Zhang, Peng Ye, Shi Feng, Daling  
652 Wang, Zhen Wang, Xinrun Wang, et al. The avengers: A simple recipe for uniting smaller  
653 language models to challenge proprietary giants. *arXiv preprint arXiv:2505.19797*, 2025e.

654  
655 Richard Zhuang, Tianhao Wu, Zhaojin Wen, Andrew Li, Jiantao Jiao, and Kannan Ramchandran.  
656 Embedllm: Learning compact representations of large language models. In *The Thirteenth In-*  
657 *ternational Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025,*  
658 *2025*.

659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A MORE EXPERIMENTAL RESULTS (ALL NEW)

### A.1 DETAILED EXPERIMENTAL SETTING AND COMPUTATIONAL COST

We trained our IrtNet for 15 epochs with a learning rate of  $1e-4$ , a weight decay of  $1e-4$ , a batch size of 2048, and a dropout rate of 0.5.

We used a single NVIDIA A100 (80GB) GPU. Training with a batch size of 2048 consumed approximately 2GB of VRAM and took a total of 15 minutes for 15 epochs. For inference, when the batch size was set equal to the number of models (112), the VRAM consumption was 1GB. Processing 3000 test samples took 20 seconds, resulting in an average time of 6.7ms per test sample.

### A.2 MODEL ROUTING

#### A.2.1 MORE DATASETS

To verify the applicability of IrtNet on more datasets, we introduce 10 datasets used by Avengers (Zhang et al., 2025e), including BBH (Suzgun et al., 2023), EmoryNLP (Byrkjeland et al., 2018), FinQA (Chen et al., 2021), HumanEval (Wang et al., 2024), K&K (Xie et al., 2024), Livecodebench (Jain et al., 2024), Math500 (Lightman et al., 2023), MBPP (Austin et al., 2021), MELD (Poria et al., 2019), and MMLUPro (Wang et al., 2024), totaling 8,550 questions. Of these, 6,840 questions (80%) were used for the training set and 1,710 questions (20%) for the test set.

In this supplementary experiment, we use more recent LLMs compared to those in the main text, specifically 20 models with 7-9B parameters. These models are: DeepHermes-3-Llama-3-8B-Preview, DeepSeek-R1-0528-Qwen3-8B, DeepSeek-R1-Distill-Qwen-7B, Fin-R1, GLM-Z1-9B-0414, Intern-S1-mini, Llama-3.1-8B-Instruct, Llama-3.1-8B-UltraMedical, Llama-3.1-Nemotron-Nano-8B-v1, MiMo-7B-RL-0530, MiniCPM4.1-8B, NVIDIA-Nemotron-Nano-9B-v2, OpenThinker3-7B, Qwen2.5-Coder-7B-Instruct, Qwen3-8B, cogito-v1-preview-llama-8B, gemma-2-9b-it, glm-4-9b-chat, granite-3.3-8b-instruct, and internlm3-8b-instruct.

Table 6: Model routing accuracy (%) comparison on Avengers datasets. **Micro** denotes the micro-averaged accuracy, calculated on the total correct predictions across all datasets. **Macro** denotes the macro-averaged accuracy, computed by first calculating the accuracy for each benchmark individually and then averaging these benchmark scores. We use **bold** to indicate the best results.

Method	Bbh	Emory	Finqa	HE	K&K	Lcb	M500	Mbpp	Meld	MP	Micro	Macro
RouterDC	65.4	<b>46.4</b>	59.1	68.0	80.6	16.7	91.3	63.5	54.6	55.7	56.9	60.1
ModelSAT	<b>90.1</b>	41.2	<b>73.7</b>	72.0	73.9	<b>72.2</b>	93.3	74.7	52.2	70.7	70.7	71.4
EmbedLLM	87.4	44.0	70.5	70.0	70.6	67.8	<b>95.3</b>	70.0	52.2	<b>71.7</b>	69.0	69.9
Avengers-Pro	88.9	42.6	<b>73.7</b>	72.0	75.4	69.7	94.7	75.4	<b>56.5</b>	69.7	71.2	71.8
IrtNet (Ours)	89.8	44.5	73.4	<b>78.0</b>	<b>81.0</b>	69.7	94.0	<b>75.8</b>	55.7	<b>71.7</b>	<b>72.1</b>	<b>73.4</b>

As shown in Table 6, our method achieves the best performance on both micro-average and macro-average accuracy across 10 benchmarks. Combining the experiments in Table 1, IrtNet achieves SOTA routing results on a total of 20 datasets, fully demonstrating the versatility of our method.

#### A.2.2 NOISY SETTING

To investigate the effect of noisy data on IrtNet, we adopt the original dataset from EmbedLLM. In this original dataset (including both the training and test sets), the same question might be simultaneously assigned multiple labels of both 0 and 1. The experimental results, shown in Table 7, indicate that under noisy data, when compared to the results in Table 1, our routing performance does not decline and still maintain the SOTA level, demonstrating that IrtNet can function effectively even in noisy environments.

Table 7: Model routing accuracy (%) comparison on Embedllm dataset with noisy labels.

Method	ASD	GPQA	GSM	Math	Logi	Med	Mmlu	Soci	PIQA	Tru	Micro	Macro
RouterDC	64.7	26.6	70.5	38.4	39.2	59.3	72.9	27.8	81.3	46.0	56.9	52.7
EmbedLLM	48.0	28	82.1	48.5	43.1	65.5	80.8	<b>35.2</b>	<b>88.1</b>	50.0	62.0	56.9
MODEL-SAT	8.1	11.0	67.9	39.7	45.1	56.5	<b>86.3</b>	32.1	86.6	33.8	55.8	46.7
Avengers-Pro	13.6	12.8	<b>89.3</b>	55.7	<b>49.0</b>	72.3	84.0	30.9	86.6	41.9	59.7	53.6
IrtNet (Ours)	<b>66.7</b>	<b>27.8</b>	<b>89.3</b>	<b>57.0</b>	47.1	<b>75.4</b>	86.2	30.9	85.8	<b>52.7</b>	<b>67.5</b>	<b>62.7</b>

### A.3 BENCHMARK PREDICTION

#### A.3.1 MORE BASELINES

In this section, we consider comparing additional benchmark prediction methods. We primarily evaluate three approaches: Sloth (Polo et al., 2024a), observational scaling laws (Ruan et al., 2024), and Random-Sampling-Learn (Zhang et al., 2025a). We focus on comparing the latter two because Sloth requires two key inputs: model parameters and the number of training tokens. While parameter counts are accessible for most open-sourced models, the specific number of training tokens is often unavailable.

We maintain the OOD experimental setup to evaluate the RMSE between the predicted and actual benchmark scores across all methods. To enable the implementation of OSL and RSL, we have to reserve an additional 5% of the models, utilizing their performance on 50 cold-start queries along with their ground-truth full-benchmark scores for training. In contrast, EmbedLLM and IrtNet do not require this additional supervision, as they rely solely on query semantics and model identifiers for prediction.

Table 8: Root mean square error (RMSE) for OOD benchmark prediction. **Avg.** denotes the average RMSE of all benchmarks.

Method	GPQA	GSM	Math	Logi	Med	Mmlu	PIQA	Tru	Avg.
OSL (Ruan et al., 2024)	0.30	0.36	<b>0.06</b>	0.21	0.28	0.45	0.51	0.36	0.32
RSL (Zhang et al., 2025a)	<b>0.23</b>	0.43	0.28	<b>0.08</b>	0.27	0.37	<b>0.35</b>	0.29	0.29
EmbedLLM (Zhuang et al., 2025)	0.26	0.20	0.09	0.11	0.04	0.25	0.40	<b>0.11</b>	0.18
IrtNet (Ours)	0.26	<b>0.19</b>	0.09	0.11	0.04	<b>0.19</b>	0.36	0.12	<b>0.16</b>

The results are presented in Table 8. It is evident that OSL and RSL exhibit significantly higher prediction errors compared to EmbedLLM and IrtNet. This disparity primarily stems from the difference in prediction granularity: OSL and RSL perform benchmark-level predictions by inferring the relationship between few-shot observations and full-benchmark scores based on historical models, whereas EmbedLLM and IrtNet achieve significantly greater precision through fine-grained, query-level prediction.

#### A.3.2 MORE RESULTS IN TABLE 3

In this subsection, we additionally present the result figures for the remaining five OOD Benchmark predictions that were not shown in Section 4.3. As shown in Figure 6, on these five datasets, the prediction results of IrtNet and EmbedLLM are essentially tied, and they exhibit the same bias, meaning they either simultaneously overestimate or underestimate the results for a specific dataset.

### A.4 ABLATION STUDY

#### A.4.1 QUERY ENCODER

We perform an ablation study on the query encoder, replacing all-mpnet-base-v2 with the relatively weaker bert-base-nli-mean-tokens (Reimers & Gurevych, 2019). As shown in Table 9, compared to the results in Table 1, all methods experience performance fluctuations. Except for ModelSAT, which sees a performance increase, all other methods show a decline in performance. This suggests

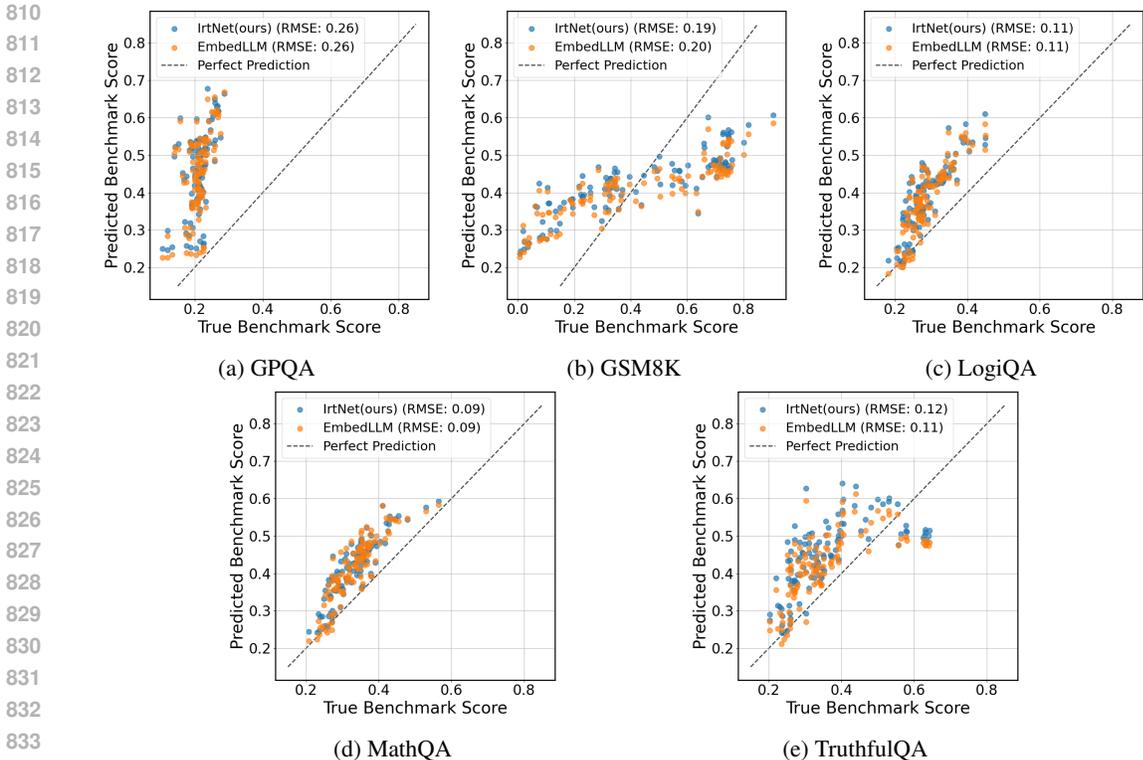


Figure 6: Predicted vs. true benchmark scores (in [0-1]) on five OOD benchmarks. The scatter plots represent the predicted LLM scores by IrtNet and EmbedLLM.

that routing methods are indeed dependent on the effectiveness of the query encoder to some extent. Overall, we still maintain strong competitiveness, achieving first place in Micro Accuracy and second in Macro Accuracy.

Table 9: Model routing accuracy (%) comparison with bert-base-nli-mean-tokens as query encoder.

Method	ASD	GPQA	GSM	Math	Logi	Med	Mmlu	Soci	PIQA	Tru	Micro	Macro
RouterDC	36.9	16.6	21.4	40.9	39.2	59.0	65.8	28.4	81.3	46.0	49.0	43.6
EmbedLLM	15.7	22.0	<b>88.4</b>	<b>54.4</b>	47.1	67.5	79.1	30.3	85.8	48.7	58.8	53.9
MODEL-SAT	<b>70.2</b>	26.6	86.6	49.0	45.1	69.5	82.3	29.6	85.1	<b>50.0</b>	64.1	<b>59.4</b>
Avengers-Pro	11.1	<b>28.6</b>	<b>88.4</b>	<b>54.4</b>	<b>49.0</b>	<b>73.2</b>	81.7	<b>30.9</b>	<b>86.6</b>	35.1	61.0	53.9
IrtNet (Ours)	61.6	25.8	86.6	53.2	41.2	70.6	<b>85.0</b>	29.6	85.8	<b>50.0</b>	<b>64.9</b>	58.9

#### A.4.2 MODEL EMBEDDING DIMENSION

In this section, we investigate the sensitivity of IrtNet to the dimension of the LLM’s compact representation. We only vary the dimension of the learned model embedding and test the performance on both the model routing and benchmark prediction tasks.

Table 10: Sensitivity of IrtNet to the dimension of the LLM Embedding. **Routing** denotes the model routing task and **Correctness** denotes the correctness prediction task.

Task	LLM Embedding Dimension						
	32	64	128	232 (Ours)	512	1024	2048
<b>Routing</b>	64.8	65.9	66.9	<u>67.4</u>	66.7	67.2	<b>67.6</b>
<b>Correctness</b>	71.8	71.9	72.0	<b>72.2</b>	<u>72.1</u>	72.0	72.0

As shown in Table 10, using 232 dimensions as the baseline (chosen for consistency with EmbedLLM), when we decrease the embedding dimension, the performance of both tasks significantly drop. This is primarily because an overly small dimension cannot capture sufficient information. Conversely, when we continue to increase the model embedding dimension, the performance barely improves. This suggests that due to limited training data, further increasing the dimension becomes redundant once a suitable dimension is sufficient to fully represent the information.

#### A.4.3 MOE EXPERT NUMBER

In this subsection, we investigate the effect of MoE sparsity/density and the influence of the number of experts. We adopted two settings: In the dense setting, we simultaneously varied the total number of experts and the number of activated experts. In the sparse setting, we kept the total number of experts at 40 and then varied the number of activated experts.

As shown in Table 11, due to our limited data size, the sparse MoE indeed suffered from insufficient training, and the model’s performance gradually improved as the number of activated experts increased. In the dense setting, the model’s performance change was largely insignificant.

Table 11: Ablation study of MoE expert number. We use **bold** to indicate the best results.

MoE Architecture	Total Experts	Active Experts	Routing	Correctness
Dense (Ours)	40	40	<b>67.4</b>	<b>72.2</b>
Sparse	40	4	62.4 $\downarrow$ 5.0	68.7 $\downarrow$ 3.5
Sparse	40	20	66.1 $\downarrow$ 1.3	71.2 $\downarrow$ 1.0
Dense	20	20	67.0 $\downarrow$ 0.4	72.0 $\downarrow$ 0.2
Dense	80	80	66.7 $\downarrow$ 0.7	72.0 $\downarrow$ 0.2

### A.5 INTERPRETABILITY ANALYSIS

#### A.5.1 UNDERSTANDING DISCRIMINATION

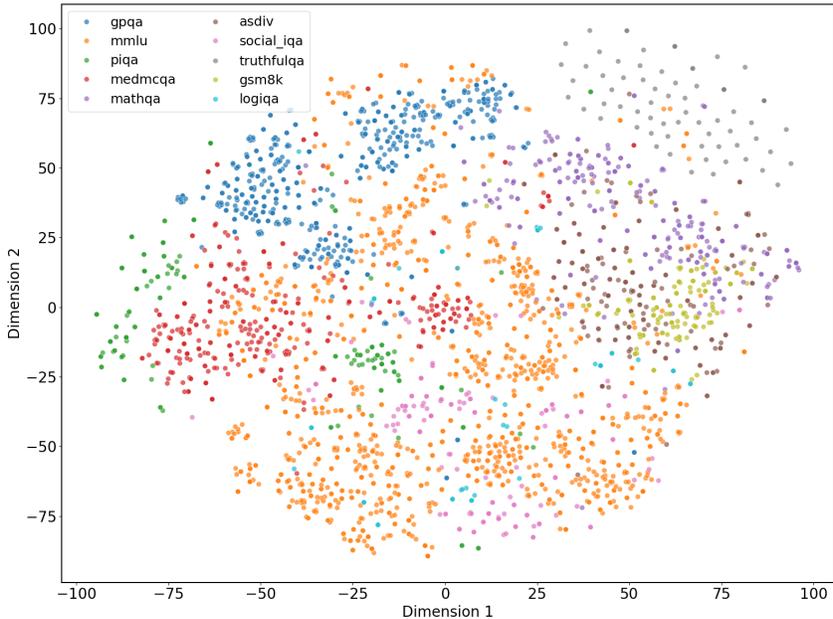


Figure 7: T-SNE visualization of raw query embeddings.

As shown in Figure 7, we perform the t-SNE visualization on the raw Sentence-BERT embeddings ( $\mathbf{v}_q$ ). While the raw  $\mathbf{v}_q$  embeddings show a general clustering tendency, the clusters are often diffuse and overlapping. In contrast, the learned discrimination vectors  $\alpha_q$  in Figure 4 exhibit significantly

tighter and clearer separation between different benchmark groups. This comparison confirms that our MoE-based network successfully acts as a structural filter, refining the generic semantic features of  $\mathbf{v}_q$  into highly distinguishable and task-specific latent requirements. This refinement is the core reason for IrtNet’s superior routing and benchmark prediction performance.

A.5.2 VALIDATING DIFFICULTY

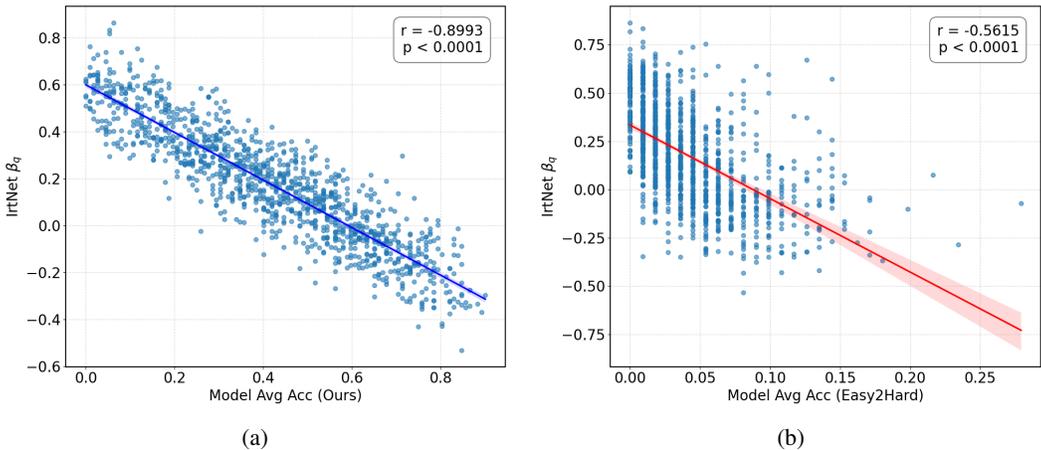


Figure 8: The correlation between the average model accuracy per question on the GSM8K dataset and the difficulty parameter  $\beta_q$ . (a) shows the correlation between the accuracy of the models used in our experiment and  $\beta_q$ , and (b) shows the correlation between the Easy2Hard-Bench model accuracy and  $\beta_q$ .

To further demonstrate the significance of the difficulty parameter learned by IrtNet, we align the GSM8K results in our experiments with the Easy2Hard-Bench (Ding et al., 2024). We calculate the correlation between the average accuracy of our model set and the learned difficulty parameter  $\beta_q$ , as well as the correlation between the average accuracy of the Easy2Hard-Bench model set on GSM8K and  $\beta_q$ .

As shown in Figure 8, the average accuracy of our model set exhibits an extremely strong negative correlation with the difficulty parameter  $\beta_q$ , with a correlation coefficient  $r = -0.8993$ . The average accuracy per question in the Easy2Hard-Bench model set also shows a moderate negative correlation with  $\beta_q$ , with a correlation coefficient  $r = -0.5615$ . We believe it is reasonable for the Easy2Hard-Bench correlation to be lower than the correlation within our own dataset because  $\beta_q$  fundamentally learns the average performance of the models in the training data, and the average performance of the Easy2Hard-Bench model set differs from that of our model set. This explains why our  $\beta_q$  correlates exceptionally well with the average performance of the models in our own dataset.

Concurrently, we also categorize the questions into five difficulty levels—Easiest, Easy, Medium, Hard, and Hardest—based on the model accuracy in Easy2Hard-Bench, and sampled representative questions for layered display. As shown in Table 12, the accuracies for the "Easy" and "Hard" questions are even reversed in our dataset: they are 0.063 and 0.018 in Easy2Hard-Bench, respectively, but 0.32 and 0.38 in ours. This is why the correlation between  $\beta_q$  and Easy2Hard-Bench accuracy is lower. In contrast,  $\beta_q$  correlates very well negatively with our own accuracy: the higher the accuracy, the lower the difficulty  $\beta_q$ .

A.6 MODEL LIST

Here is an exhaustive list of models contained in Section 4:

Table 13: The exhaustive list of the 112 models used in Section 4.

Model Name	Model Name
meta-llama/LlamaGuard-7b	meta-llama/Llama-2-13b-chat-hf

Table 13: The exhaustive list of the 112 models used in Section 4.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Model Name	Model Name
01-ai/Yi-34B-Chat	meta-llama/Llama-2-70b-chat-hf
WizardLM/WizardLM-70B-V1.0	allenai/tulu-2-dpo-70b
Imsys/vicuna-13b-v1.5	Imsys/vicuna-33b-v1.3
Qwen/Qwen-14B-Chat	upstage/SOLAR-10.7B-Instruct-v1.0
openchat/openchat-3.5-0106	openchat/openchat-3.5
berkeley-nest/Starling-LM-7B-alpha	HuggingFaceH4/zephyr-7b-beta
TheBloke/tulu-30B-fp16	mistralai/Mistral-7B-Instruct-v0.1
tiuae/falcon-40b-instruct	Imsys/vicuna-13b-v1.5-16k
codellama/CodeLlama-34b-Instruct-hf	TheBloke/WizardLM-13B-V1.2-GGUF
Imsys/vicuna-7b-v1.5	NousResearch/Nous-Hermes-13b
project-baize/baize-v2-13b	Imsys/vicuna-7b-v1.5-16k
mosaicml/mpt-30b-instruct	meta-llama/Llama-2-7b-chat-hf
TheBloke/koala-13B-HF	nomic-ai/gpt4all-13b-snoozy
h2oai/h2ogpt-gm-oasst1-en-2048-open-llama-13b	mosaicml/mpt-7b-chat
databricks/dolly-v2-12b	stabilityai/stablelm-tuned-alpha-7b
OpenAssistant/oasst-sft-4-pythia-12b-epoch-3.5	deepseek-ai/deepseek-llm-67b-chat
NousResearch/Nous-Hermes-2-Yi-34B	CausalLM/34b-beta
SUSTech/SUS-Chat-34B	SUSTech/SUS-Chat-72B
Qwen/Qwen-72B	Intel/neural-chat-7b-v3-3
ibivibiv/alpaca-dragon-72b-v1	JaeyeonKang/CCK-Asura-v1
ConvexAI/Luminex-34B-v0.2	ConvexAI/Luminex-34B-v0.1
CorticalStack/pastiche-crown-clown-7b-dare-dpo	ogno-monarch-jaskier-merge-7b-OH-PREF-DPO
bardsai/jaskier-7b-dpo-v5.6	FelixChao/Scorpio-7B
dfurman/HermesBagel-34B-v0.1	kevin009/llamaRAGdrama
sail/Sailor-7B	AiMavenAi/Prometheus-1.3
Q-bert/Optimus-7B	cognitivecomputations/yayi2-30b-llama
zhengr/MixTAO-7Bx2-MoE-v8.1	fbllgit/UNA-SimpleSmaug-34b-v1beta
mistralai/Mixtral-8x7B-Instruct-v0.1	microsoft/Orcas-2-13b
EleutherAI/pythia-12b	cloudyu/Mixtral-11Bx2-MoE-19B
rishiraj/CatPPT-base	Deci/DeciLM-7B
microsoft/phi-2	scb10x/typhoon-7b
01-ai/Yi-6B-200K	01-ai/Yi-6B
TigerResearch/tigerbot-13b-base	augmxnt/shisa-base-7b-v1
microsoft/phi-1.5	golaxy/gowizardlm
bigscience/bloom-7b1	mlabonne/AlphaMonarch-7B
Cultrix/NeuralTrix-bf16	shadowml/MBeagleX-7B
yam-peleg/Experiment26-7B	deepseek-ai/deepseek-math-7b-instruct
meta-math/MetaMath-Mistral-7B	kyujinpy/Sakura-SOLRCA-Math-Instruct-DPO-v1
FelixChao/llama2-13b-math1.2	Plaban81/Moe-4x7b-math-reason-code
MazyarPanahi/WizardLM-Math-70B-v0.1	abhishek/zephyr-beta-math
meta-math/MetaMath-Llemma-7B	EleutherAI/llemma-34b
EleutherAI/llemma-7b	FelixChao/vicuna-7B-physics
Harshvir/Llama-2-7B-physics	FelixChao/vicuna-7B-chemical
BioMistral/BioMistral-7B	BioMistral/BioMistral-7B-DARE
PharMolix/BioMedGPT-LM-7B	Biomimicry-AI/ANIMA-Nectar-v2
codellama/CodeLlama-7b-hf	codellama/CodeLlama-13b-Instruct-hf
deepseek-ai/deepseek-coder-1.3b-base	deepseek-ai/deepseek-coder-6.7b-instruct
OpenBuddy/openbuddy-codellama2-34b-v11.1.1-bf16	TheBloke/CodeLlama-70B-Instruct-AWQ
AdaptLLM/medicine-chat	AdaptLLM/medicine-LLM
AdaptLLM/medicine-LLM-13B	Writer/palmyra-med-20b
SciPhi/SciPhi-Self-RAG-Mistral-7B-32k	Neko-Institute-of-Science/metharme-7b
Neko-Institute-of-Science/pygmalion-7b	SciPhi/SciPhi-Mistral-7B-32k
shleeeee/mistral-ko-tech-science-v1	codefuse-ai/CodeFuse-DeepSeek-33B
WizardLM/WizardCoder-Python-34B-V1.0	bigcode/octocoder
meta-llama/Meta-Llama-3-8B	meta-llama/Meta-Llama-3-8B-Instruct
meta-llama/Meta-Llama-3-70B	meta-llama/Meta-Llama-3-70B-Instruct

Table 13: The exhaustive list of the 112 models used in Section 4.

Model Name	Model Name
meta-llama/Meta-Llama-Guard-2-8B	Qwen/Qwen1.5-32B-Chat
Qwen/Qwen1.5-4B-Chat	Qwen/Qwen1.5-0.5B-Chat
Qwen/Qwen1.5-7B-Chat	Nexusflow/Starling-LM-7B-beta
google/gemma-7b-it	google/gemma-2b-it

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

Table 12: Case study of question accuracy and difficulty parameter  $\beta_q$ .

<b>E2H Difficulty</b>	<b>Question</b>	<b>E2H Acc</b>	<b>Our Acc</b>	<b>IrtNet <math>\beta_q</math></b>
Easiest	A raspberry bush has 6 clusters of 20 fruit each and 67 individual fruit scattered across the bush. How many raspberries are there total?	0.2793	0.7054	-0.0716
Easy	It's strawberry-picking time on Grandma Concetta's farm. Tony can pick 6 quarts of strawberries per hour, while Bobby picks one less quart of strawberries per hour than Tony. Kathy can pick twice as many strawberries per hour as Bobby, and Ricky picks two fewer quarts of strawberries per hour than does Kathy. In total, how many quarts of strawberries can Tony, Bobby, Ricky, and Kathy pick per hour on Grandma Concetta's farm?	0.0631	0.3214	0.3380
Medium	Laurel's friend gave her 24 baby outfits that her child no longer needed. At her baby shower, Laurel received twice the amount of new baby outfits. Then, Laurel's mom gifted her with another 15 baby outfits. How many outfits does she have for her baby?	0.0360	0.4643	0.1699
Hard	Oscar has 24 lollipops and eats 2 on his way to school. He passes 14 out to his friends. He buys twice as many lollipops on his way home as he gave to his friends. He eats 3 more that night and 2 more in the morning. How many lollipops does Oscar have?	0.0180	0.3839	0.2611
Hardest	For a New Year's resolution, Andy wants to lose 30 lbs. by his birthday, which is July 19th. Today is December 31st. If Andy needs to burn 3500 calories to lose a pound, how much of a calorie deficit (net amount of calories burned vs. calories consumed) does he need each day to reach his goal?	0.0000	0.0268	0.6406