# Complexity-aware fine-tuning

**Andrey Goncharov[1]** ⓘ**, Daniil Vyazhev[1]** ⓘ**, Petr Sychev[1],**
**Edvard Khalafyan, Alexey Zaytsev[1]** ⓘ**,**
[1]Applied AI Institute

## Abstract

General-purpose Large Language Models (LLMs) are frequently fine-tuned through supervised fine-tuning (SFT) to enhance performance in specific domains. Better results can be achieved by distilling the chain-of-thought of a larger model at the cost of numerous expensive calls and a much greater amount of data. We propose a novel blueprint for efficient fine-tuning that uses reasoning only for complex data identified by entropy. Specifically, across two small open models ($\approx 3B$) we split the training data into complexity categories by a single token answer entropy (ROC AUC $0.73$), fine-tune large language models (LLMs) via SFT and distillation, and show that our pipeline significantly outperforms the standard SFT approach ($0.58$ vs $0.45$ average accuracy) and outperforms the distillation approach ($0.58$ vs $0.56$ average accuracy) while using $81\%$ less data. We publish our code[1] and data[2] to facilitate further research in this direction.

## 1 Introduction

General-purpose LLMs excel across diverse tasks, but deploying them is often impractical under constraints on compute, latency or cost. These realities motivate compact, domain-adapted models that can deliver competitive or superior performance. Growing evidence shows that carefully tuned smaller models can match or outperform larger open models in mathematics (Yang et al., 2024b), medicine (Wu et al., 2025), chemistry (Yu et al., 2024). Thus, smaller domain-specific LLMs appear a compelling choice, especially under resource constraints.

A standard way for training is domain-adaptive fine-tuning: starting from a general-purpose LLM and training it on in-domain data (Parthasarathy

[1]https://github.com/LabARSS/
complexity-aware-fine-tuning/tree/eacl2026
[2]https://github.com/LabARSS/
complexity-aware-fine-tuning/tree/eacl2026?
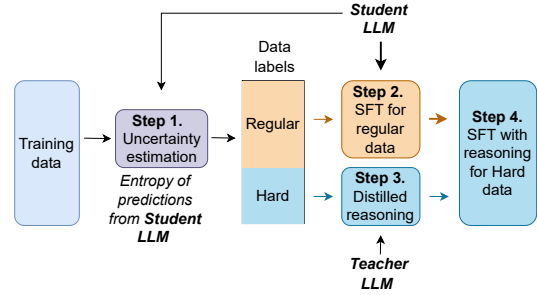tab=readme-ov-file#data



Figure 1: Complexity-aware fine-tuning scheme for a student LLM: we identify complexity of questions via uncertainty estimation of a model (Step 1), then for questions of regular complexity we apply direct SFT (Step 2), while for hard questions we include reasoning from a teacher LLM (Step 3) and complete SFT using reasoning-enriched hard data (Step 4).

et al., 2024) with optional distillation from a larger model. Recent work has mainly concentrated on optimizing training mechanics (objectives, schedulers, and regularizers), whereas comparatively little attention has been paid to data curation. Curriculum-style approaches have also been explored (Kim and Lee, 2024; Shi et al., 2025) with increasing complexity of data during training, but with limited gains in final accuracy and efficiency, suggesting that more principled strategies are needed.

To address this gap, we propose a fully automated pipeline, presented in Figure 1, that focus on both accuracy and efficiency of the fine-tuning procedure. It consists of two steps: (1) a split of the data for fine-tuning by the complexity of the task to regular and hard questions and (2) training a model via supervised fine-tuning (SFT) followed by another round of SFT that utilizes distillation on a chain-of-thought from a larger model (Li et al., 2023; Hsieh et al., 2023). We confirm the effectiveness of the framework by fine-tuning two open models: Qwen2.5-3B (Yang et al., 2024a) and Phi-4-Mini (Microsoft et al., 2025), on the multi-

ple choice question answering dataset MMLU-Pro (Wang et al., 2024).

For both steps of the pipeline, we provide sensitivity studies to the design choices made. In step one, we consider different options for complexity estimation based on a model's confidence in its answer. Our detailed study considers a wide range of methods, including model-as-judge (MASJ), entropy-based aggregation, and calculation methods in line with findings from Fadeeva et al. (2023) with a top performing one being an answer entropy. In step two, we explore training on different complexity subsets of the distilled chain-of-thought — concluding that for regular data standard SFT is enough, while for hard data we benefit reasoning-based distillation. This observation constrains the reasoning to hard data only, making our approach token-efficient.

The specific contributions are the following:

- A training pipeline to fine-tune LLMs via SFT and distillation for regular and hard data correspondingly. Our procedure outperforms the standard SFT approach with accuracies 0.52 and 0.64 over two LLMs for MMLU-Pro c.t. 0.39 and 0.51 for SFT and 0.50 and 0.63 for distillation baselines.

- An unsupervised method to split a multiple-choice question answering dataset by complexity based on the token-wise entropy of the response with ROC AUC 0.73. Alternative complexity estimators rooted in MASJ and aggregated entropy, augmented with reasoning results and various aggregation techniques yield inferior results.

- Open-source standardized datasets to facilitate further development of uncertainty estimation and calibration methods: with and without chain-of-thought, with token probability distribution at each step provided, as well as additional scores.

## 2 Related works

**Data complexity-aware learning.** Curriculum learning has been explored to improve LLM fine-tuning by ordering training examples from easy to hard. Kim and Lee (2024) propose sorting fine-tuning data by difficulty metrics (e.g. prompt length, model attention scores, and initial loss) so that the model learns on simpler prompts before complex ones. They found that this curriculum strategy yielded slightly higher accuracy than random shuffling, with ordering by an attention-based criterion performing best. This approach is attractive because it boosts performance without adding more data or parameters. However, the gains were modest, and defining difficulty automatically can be tricky - their method requires measuring things like loss or attention per example.

Another strategy is filtering training data for quality. A notable example is LIMA (Zhou et al., 2023a), which shows that a large pre-trained model can be fine-tuned on just a small, high-quality subset of data. They fine-tuned a 65B Llama model on only 1000 carefully curated prompt-response pairs (chosen for diversity and clarity) without any reinforcement learning. Despite the tiny dataset, the resulting model performed remarkably well, learning to handle complex queries and even generalizing to tasks not seen in training. In a human evaluation, LIMA's answers were preferred over GPT-4's in 0.43 of cases. This "less is more" result suggests that much of an LLM's ability comes from pre-training, and fine-tuning needs only a small amount of exemplary data to unlock it. However, LIMA relied on a large base model and manual data curation. The approach may not scale down to smaller models and requires human intervention.

Another notable example of curated data selection is the SmallToLarge (S2L) method by Yang et al. (2024c), which leverages training trajectories from small models to guide the data selection for larger models. This way, the large LLM is trained on a diverse yet compact dataset covering different difficulty levels. S2L showed impressive results: for a math word problem dataset, they achieved the same accuracy using only 11% of the data, and even outperformed other selection methods by 4.7% on average across several benchmarks. The strength of this approach is that it makes complexity-based data filtering automated and cheap. One caveat is the extra step of training a smaller model and clustering. The approach is mostly tested on specialized domains (math problems, clinical text summarization), so its generality to all types of tasks needs further validation. Additionally, it requires a large amount of data to make a filtered subset.

**Complexity estimation.** Fadeeva et al. (2023) provides an in-depth comparison of multiple black-box and white-box methods of complexity estimation. They present promising results for white-box methods rooted in sequence probability and entropy

aggregation. Sychev et al. (2025) focus specifically on entropy-based aggregations augmented with model-as-judge (MASJ) categorization by a reasoning score. They confirm that LLM's token-level entropy of the output is a good predictor of question difficulty, especially in knowledge-based domains. They also introduce MASJ reasoning score to estimate the question complexity. However, the authors use these metrics only to analyze model behavior. They do not integrate it into a practical data aggregation or fine-tuning workflow.

**Research gap.** The question if we can combine novel complexity estimation techniques based on entropy aggregates with adaptive learning methods remains open. It is also unclear how we can use the insights for the reasoning estimates to dynamically change our learning approach. This paper aims to address these gaps.

## 3 Methods

### 3.1 Training pipeline

We propose the complexity-aware fine-tuning pipeline detailed in Figure 1 with the following major stages: complexity estimation, data aggregation, fine-tuning.

**Complexity estimation.** We adopt the entropy of the answer token in the response as our primary complexity metric. We prompt the model to pick the correct option directly, without producing a chain-of-thought. See Section 3.2.3 for details.

**Data aggregation.** To aggregate the data into groups by complexity (regular, hard), we evenly divide the dataset into three parts, ordering the entries by entropy of the response. A group with lower entropy values is categorized as easy and with the higher values as hard, as Figure 2 depicts.

**Fine-tuning.** We propose to apply different fine-tuning strategies according to the complexity of the data.

For the regular groups, we use vanilla SFT (Howard and Ruder, 2018; Raffel et al., 2023), an established and robust practice. It involves fine-tuning a pretrained LM on labeled examples using a standard supervised objective, specifically like cross-entropy. Specific prompts are provided in Table 4 in Appendix.

As to the hard group, we hypothesize that as hard questions require multiple logical steps and multiple attempts to learn core facts. So, the standard
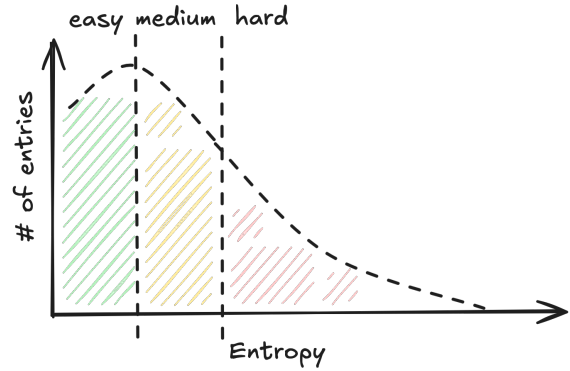


Figure 2: Data aggregation: we can further split the data to various complexity chunks

SFT is suboptimal. We propose to elicit a chain-of-thought and allow the model to incrementally build the answer step-by-step as suggested by Wei et al. (2023).

In this work, we apply the distillation technique, which involves training a smaller student model on the chain-of-thought of a larger LLM. It is a well-known practice supported by (Hsieh et al., 2023). To create the distillation training samples, we prompt a large LLM to answer the multiple-choice question and produce a chain-of-thought in the process. Next, the whole response is attached to the dataset and used to train the smaller model.

### 3.2 Complexity estimation approaches

To find the most suitable metric for the training pipeline, we analyze the performance of the following techniques: MASJ reasoning score, MASJ education level, a single token answer entropy, a chain-of-thought answer entropy, a chain-of-thought aggregated response entropy. Used prompts are available in Appendix A.1.

#### 3.2.1 MASJ reasoning score

As one of the expert-like metrics, we ask a large LLM to estimate the number of logical steps required to answer the question. The hypothesis is that the questions that require more reasoning should be harder for the model to answer.

To collect the MASJ-based reasoning score, we go over the multiple choice question answering dataset and query a large auxiliary LLM for the estimate. We prompt the model to provide the number of logical steps required to answer the question: low, medium, high. Next, we query the large LLM again to estimate the quality of the previous assessment from 1 to 10 following the practice introduced

in MT-Bench by Zheng et al. (2023). It allows us to filter out low quality scores by keeping only the ones rates above or equal to 9.

### 3.2.2 MASJ education level

As the other expert-like metrics, we ask a large LLM to estimate the required level of education to answer the question correctly. It is a natural human-like value used in other datasets (Rein et al., 2023; Lu et al., 2022).

We follow the same procedure as for the MASJ reasoning score, but use a different prompt.

### 3.2.3 Answer entropy

**Single token answer entropy.** In a similar fashion to that proposed by (Kadavath et al., 2022), we calculate the entropy of the answer token in the response. The assumption is that the response uncertainty is a natural predictor of the question complexity. We prompt the model to answer the question directly (as a single token) and calculate token-wise entropy of the response as follows:

$$h = -\sum_{i=1}^{n} p_i \log p_i,$$

where $p_i$ is the probability of a single token and $n$ is the vocabulary size. This is the main method used in our pipeline.

Additionally, similarly to (Zhou et al., 2023b), we examine the performance when we allow the model to explicitly say "I do not know" (IDK).

**Chain-of-thought answer entropy.** With the same assumption as for the single token entropy, we analyze the entropy of the answer token, but change the prompt to elicit a chain-of-thought type of response. The hypothesis is that, through the chain-of-thought process, the LLM can incrementally accumulate entropy, resulting in a better separation of certain and uncertain answers.

### 3.2.4 Chain-of-thought aggregated response entropy.

Building upon the single-token entropy approach, we investigate more sophisticated methods in line with Fadeeva et al. (2023) for complexity estimation by analyzing the entire chain-of-thought (CoT) response. While the answer token entropy provides a localized measure of uncertainty, aggregating entropy across the complete reasoning process potentially offers a more comprehensive complexity assessment.

We evaluate ten distinct aggregation methods applied to CoT responses, comparing their effectiveness through ROC AUC metrics across multiple models, with and without "I don't know" option in the prompt. We consider word aggregation, sequence aggregation, and probability-based methods. See Appendix B for details.

## 4 Results

### 4.1 Experimental setup

**Dataset** We conduct all experiments on the multiple choice question answering dataset MMLU-Pro (Wang et al., 2024), widely adopted by the community as one of the main performance benchmarks. It spans 14 domains, offering a broad selection of questions of varying complexities. Each question has approximately ten options, with a single correct one, which eliminates ambiguity in evaluation.

**LLMs** Regarding the models [3], we utilize a variety of open models to analyze how the trend changes with model size. At the same time, we focus on smaller models for fine-tuning to make our results reproducible and more relevant to the distillation scenario.

We apply our pipeline to two models: Qwen2.5-3B and Phi-4-Mini. For them, we measure single token and chain-of-thought entropy, collect other response metadata, fine-tune the models, and evaluate the overall pipeline.

Bigger models are used to further extend the entropy aggregate and MASJ analysis: Mistral 24B, Phi-4, Mistral 123B (Mistral, 2024). For the chain-of-thought distillation we utilize an ensemble of models: DeepSeek-V3-0324 (DeepSeek-AI, 2024), Qwen 3 235B, Llama 4 Maverick (Meta, 2025).

### 4.1.1 Our method and baselines

In experiments, we performed training for 10 or 20 epochs, see more information below. In case of the reasoning used during training, we split training to two equal halves for the first and the second round of fine-tuning.

**Pipeline (ours)** For easy and medium groups, we perform SFT for five epochs. For the hard group, we apply learning from a distilled chain-of-thought of a larger model for another five epochs.

---

[3] All models and datasets are published under permissive licenses that allow them to be used for research purposes

**Alternative** As the alternative approach, we perform SFT for five epochs for the hard group. Next, for easy and medium groups, we apply learning from a distilled chain-of-thought of a larger model for another five epochs.

**SFT** As our first baseline, we train the model via SFT without the data split for 10 epochs.

**Curriculum** For the second baseline, we train the model via curriculum-based SFT: 3 epochs on easy data, 3 epochs on medium data, and 4 epochs on hard data.

**Distillation** As an idealistic target, we tuned the LLM via distillation without the data split for 10 epochs.

### 4.1.2 Technical details

Based on ROC AUC results and positive performance of SFT for medium and easy questions (see 4.3), we take single token entropy as our complexity metric for the pipeline described in Section 3.1.

We randomly split the data into train, validation, and test with a ratio of 70:10:20. Next, in each chunk, we balance the data so that the number of entries in each complexity group is equal. Since there are fewer hard questions, we filter out medium and easy ones to match the size of the hard group.

### 4.2 Main results

Table 1 and Figures 3a 3b show the results. We see that the proposed training scheme results in a significant improvement over SFT, curriculum-based SFT and an alternative training scheme that uses distillation for only easy and medium questions. Qwen 3B achieves an accuracy of 0.50 compared to 0.47 and 0.42 for alternative and baseline, respectively, while Phi-4-mini gets to 0.60 compared to 0.58 and 0.46. At the same time, we see that the pipeline provides comparable performance to the distillation with significantly less data in number of tokens processed. Qwen 3B achieves an accuracy of 0.50 compared to 0.49 saving 79% tokens, while Phi-4-mini gets to 0.60 compared to 0.63 saving 82% tokens. We can also detect an uptrend in the pipeline training while the full distillation approach exhibits the signs of a plateau.

To confirm the plateau trend, we performed an extended set of experiments for 20 epochs with the same split of training epochs. Table 2 and Figures 4a 4b show the results. We find an even stronger

confirmation of the previous findings with the proposed scheme outperforming all baselines. Qwen 3B achieves an accuracy of 0.52 compared to 0.50 and 0.39 for distillation and baseline, respectively, while Phi-4-mini gets to 0.64 compared to 0.63 and 0.51. The same data savings of 79% and 82% apply.

| Method | Qwen 3B | Phi4-mini |
|---|---|---|
| SFT | 0.42 / 29k | 0.55 / 27k |
| Curriculum | 0.45 / 29k | 0.54 / 27k |
| Distillation | 0.49 / 19.72 k | **0.63** / 15.15k |
| Alternative | 0.47 / 5.88k | 0.58 / 4.91k |
| Pipeline (ours) | **0.50 / 3.99k** | 0.60 / **2.67k** |

Table 1: Accuracy / tokens processed for complexity-aware fine-tuning pipelines after 10 epochs

| Method | Qwen 3B | Phi4-mini |
|---|---|---|
| SFT | 0.39 / 59k | 0.51 / 54k |
| Distillation | 0.50 / 39.45k | 0.63 / 30.30k |
| Pipeline (ours) | **0.52 / 7.98k** | **0.64 / 5.35k** |

Table 2: Accuracy / tokens processed for complexity-aware fine-tuning pipelines after 20 epochs

### 4.3 Sensitivity study

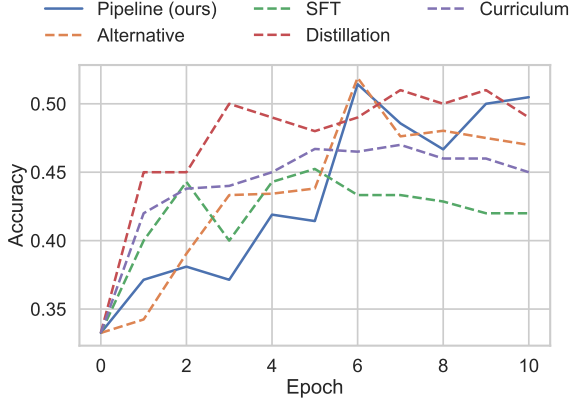#### 4.3.1 Alternative complexity estimates

To analyze the appropriate fine-tuning methods for each complexity band we start with performing the standard SFT across each group.

The same data splitting procedure as described in Section 4.2 is applied for MASJ reasoning score and single-token entropy as complexity metrics.
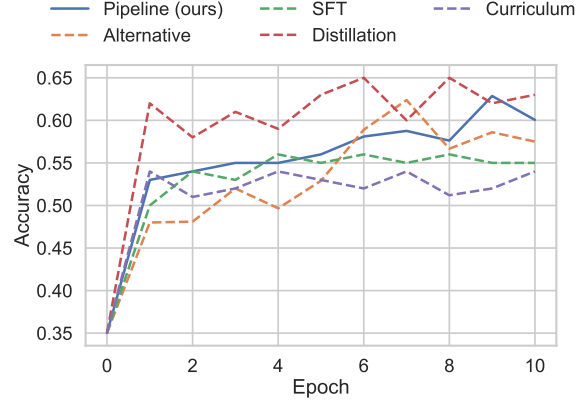
Figures 5d 5c show the results of SFT for each group split by the single-token entropy. For Phi-4-mini, we see that medium and easy questions outperform hard ones for the first 10 epochs. Shortly after, performance starts to decline for all groups. For Qwen 3B, we do not see a significant difference between the groups. Moreover, the performance plateaus after five epochs.

Figures 5b and 5a show the results of SFT for each group split by the MASJ reasoning score. We do not see a strong difference in performance between the groups. In combination with questionable ROC AUC scores, provided in Table 9, it makes MASJ reasoning score a less favorable metric for further experiments.
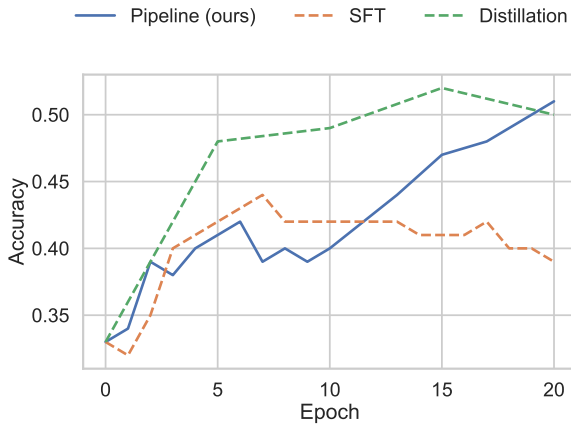
Figures 6a and 6b show pipeline training results with cross-entropy as a complexity metric com-
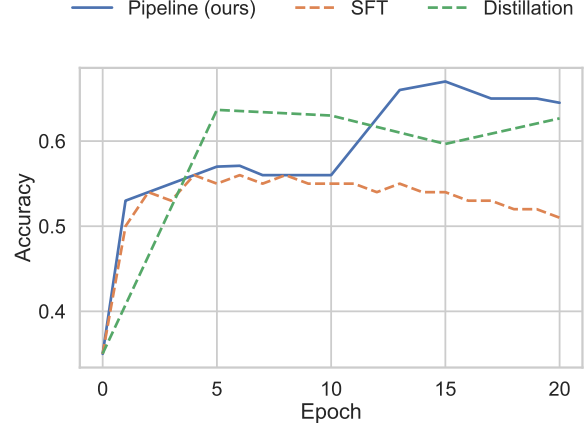
(a) Accuracy for fine-tuning pipelines after 10 epochs for Qwen 3B



(b) Accuracy for fine-tuning pipelines after 10 epochs for Phi-4-mini



(a) Accuracy for fine-tuning pipelines after 20 epochs for Qwen 3B



(b) Accuracy for fine-tuning pipelines after 20 epochs for Phi-4-mini

pared to entropy. We observe comparable results for Qwen 3B, while Phi4-mini demonstrates superior performance with entropy. The difference in performance could be attributed to entropy capturing the information across the entire distribution instead of a single correct token.

### 4.3.2 Curriculum learning fragility

Multi-stage pipelines such as curriculum learning show extreme sensitivity to hyperparameters such as a number of training epochs, learning rates and etc.

Figure 7 reveal the fragility of the original curriculum learning approach based on SFT. In this experiment we train a model on the data with increasing complexity for 10 and 20 epochs with different splits:

- 10 epochs: 3 epochs on easy data, 3 epochs on medium data, and 4 epochs on hard data.

- 20 epochs: 5 epochs on easy data, 5 epochs on medium data, and 10 epochs on hard data.

We see how Qwen 3B plateaus at a lower accuracy with more training epochs, while Phi4-mini quickly overfits on the data and its performance degrades.

### 4.3.3 Quality of complexity estimation

We consider three families of uncertainty estimation methods: MASJ, single-token entropy-based, and entropy-based augmented with a chain-of-thought. MASJ results, as they are inferior to others, are provided in Appendix C.1.

**Single token and chain-of-thought answer entropy** Table 3 presents the main ROC AUC values for single token entropy response and for the various aggregates of the chain-of-thought type of response. IDK responses and results with invalid formatting are excluded from the calculations.

We see there that the single token response shows the best results. In-depth analysis of ROC AUC and accuracy scores with splits across domains and larger selection of models can be found

| Method | Qwen 3B | Phi4-mini |
|---|---|---|
| Without split | **0.72** / 0.70 | 0.72 / **0.74** |
| Education level | | |
| High school and easier | 0.73 / 0.72 | 0.76 / 0.75 |
| Undergraduate | 0.73 / 0.71 | 0.72 / 0.77 |
| Graduate | 0.66 / 0.65 | 0.64 / 0.68 |
| Postgraduate | 0.63 / 0.52 | 0.64 / 0.63 |
| MASJ reasoning score | | |
| Low | 0.72 / 0.71 | 0.78 / 0.79 |
| Medium | 0.72 / 0.70 | 0.70 / 0.72 |
| High | 0.64 / 0.62 | 0.59 / 0.58 |

(a) ROC AUC values for single-token response

| Method | Qwen 3B | Phi4-mini |
|---|---|---|
| Answer Entropy (AE) | 0.68 / 0.67 | 0.61 / 0.58 |
| COT Mean | 0.59 / 0.58 | 0.59 / 0.63 |
| COT Max | 0.63 / 0.61 | 0.6 / 0.65 |
| Seq Mean | 0.6 / 0.59 | 0.6 / 0.62 |
| Seq Max-Mean | 0.59 / 0.58 | 0.59 / 0.61 |
| Seq Mean-Max | 0.62 / 0.6 | 0.59 / 0.62 |
| Marg Diff Mean | 0.58 / 0.57 | 0.58 / 0.61 |
| Top-2 Diff | 0.51 / 0.5 | 0.5 / 0.51 |
| COT Max - AE | 0.54 / 0.53 | 0.51 / 0.57 |
| COT Max + AE | 0.7 / 0.69 | 0.62 / 0.62 |

(b) ROC AUC values for CoT response

Table 3: ROC AUC values for single-token and CoT responses
Second result is for the alternative prompt to allow model answer "I do not know"



(a) SFT by MASJ reasoning score (Qwen 3B)

(b) SFT by MASJ reasoning score (Phi-4-mini)

(c) SFT by single token entropy (Qwen 3B)

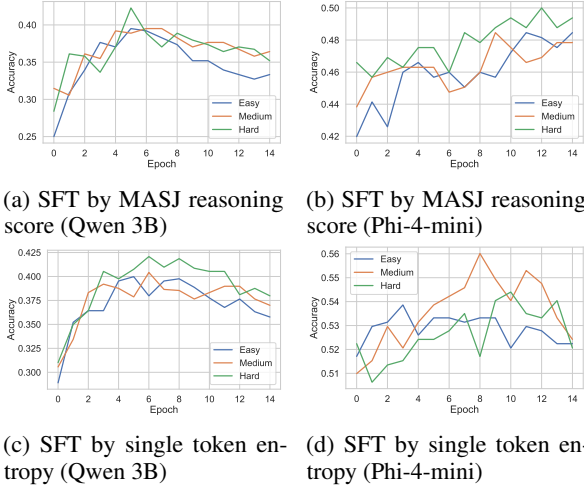(d) SFT by single token entropy (Phi-4-mini)

Figure 5: SFT quality dynamics during training with split by complexity estimates provided by the MASJ reasoning score and the single token entropy across Phi-4-mini and Qwen 3B models.



(a) Pipeline performance with cross-entropy (Qwen 3B)

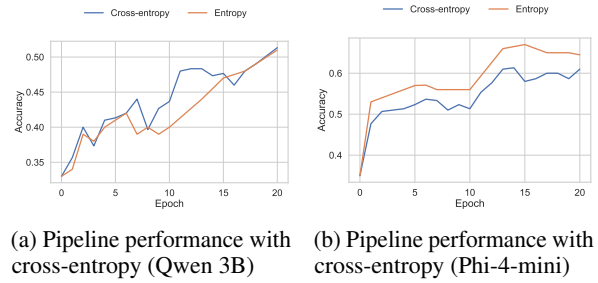(b) Pipeline performance with cross-entropy (Phi-4-mini)

Figure 6: Pipeline complexity metric performance comparison for entropy and cross-entropy across Phi-4-mini and Qwen 3B models.



Figure 7: Curriculum learning accuracy dynamics for different models for Qwen 3B (left) and Phi-4-mini (right)

in tables 7 and 8 in the Appendix.

Key observations:

- Single token response provides the best ROC AUC score. At the same time, 'IDK' responses do not consistently affect ROC AUC for all models.

- Accuracy tends to be slightly higher when an LLM can answer 'IDK'.

- Chain-of-thought responses tend to provide higher accuracy, but lower ROC AUC scores, which makes them less suitable for complexity estimation.

- Maximum-based measures (COT Max, Seq Mean Max) consistently outperform mean-based approaches (COT Mean, Seq Max Mean and Seq Mean Mean), suggesting peak uncertainty moments may better indicate question difficulty than average uncertainty.

- The poor close-to-random performance of the difference in top entropies suggests that modern LLMs maintain relatively stable reasoning to outliers.

- Sequence-based methods did not show good improvements over basic aggregation, indicating that modeling the reasoning structure provides marginal benefits.

## 5 Conclusion and discussion

This paper introduces a complexity-aware fine-tuning pipeline that measures the model response uncertainty using the entropy of its own predicted answer and then trains it on the resulting easy, medium, and hard splits with different tactics.

Using the entropy-based data split, we find that different complexity scores require different training approaches. Standard supervised fine-tuning (SFT) is enough for the easy and medium bands, but lags on the hard band. For the hard questions, adding a distilled chain-of-thought from a large LLM unlocks further gains. Our pipeline achieves accuracies of 0.52/0.64 vs. 0.39/0.51 for Qwen 3B/Phi-4-mini, while using $81\%$ less data.

The pipeline is fully automated and can be included in other fine-tuning workflows. It suggests that curriculum ideas still matter for today's LLMs: letting the model focus on what it can already solve directly, while giving extra guidance only where it struggles, yields a better allocation of limited model capacity.

We confirm that entropy works as a difficulty estimation. Single-token answer entropy reaches ROC AUC values up to $0.8$, clearly beating MASJ-based estimates of $0.57$. This confirms that a model's own confidence is a reliable, automatic proxy for question difficulty.

### Limitations

- Proposed pipeline is tested only on MMLU-Pro and small models. Results may change for other question answering datasets, open-ended tasks, other domains, or larger LLMs.

- In low-resource settings teacher may be unavailable or imperfect, which reduces the benefit of learning from a distilled chain-of-thought. Additionally, we did not explore how well the approach generalizes to other reasoning-promoting techniques.

- Low entropy can still correspond to hallucinations, which leads to imperfect identification of the question complexity.

- We split data into 3 equal parts and did not explore other possible boundaries.

- We did not conduct an extensive ablation study which might reveal that our approach does not

suggest the best possible combination or sequence of training within the current framework. It remains an area for further research.

## References

DeepSeek-AI. 2024. Deepseek-v3 technical report. *Preprint*, arXiv:2412.19437.

Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, Timothy Baldwin, and Artem Shelmanov. 2023. Lm-polygraph: Uncertainty estimation for language models. *Preprint*, arXiv:2311.07383.

Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. *Preprint*, arXiv:1801.06146.

Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. *Preprint*, arXiv:2305.02301.

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, and 1 others. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.

Jisu Kim and Juhwan Lee. 2024. Strategic data ordering: Enhancing large language model performance through curriculum learning. *arXiv preprint arXiv:2405.07490*.

Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. 2023. Symbolic chain-of-thought distillation: Small models can also "think" step-by-step. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2665–2679, Toronto, Canada. Association for Computational Linguistics.

Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.

Meta. 2025. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation.

Microsoft, :, Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benhaim, Martin Cai, Vishrav Chaudhary, Congcong Chen, Dong Chen, Dongdong Chen, Junkun Chen, Weizhu Chen, Yen-Chun

Chen, Yi ling Chen, Qi Dai, and 57 others. 2025. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *Preprint*, arXiv:2503.01743.

Mistral. 2024. Mistral-large-instruct-2411.

Venkatesh Balavadhani Parthasarathy, Ahtsham Zafar, Aafaq Khan, and Arsalan Shahid. 2024. The ultimate guide to fine-tuning llms from basics to breakthroughs: An exhaustive review of technologies, research, best practices, applied research challenges and opportunities. *Preprint*, arXiv:2408.13296.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. Exploring the limits of transfer learning with a unified text-to-text transformer. *Preprint*, arXiv:1910.10683.

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. *Preprint*, arXiv:2311.12022.

Taiwei Shi, Yiyang Wu, Linxin Song, Tianyi Zhou, and Jieyu Zhao. 2025. Efficient reinforcement fine-tuning via adaptive curriculum learning. *Preprint*, arXiv:2504.05520.

Petr Sychev, Andrey Goncharov, Daniil Vyazhev, Edvard Khalafyan, and Alexey Zaytsev. 2025. When an llm is apprehensive about its answers – and when its uncertainty is justified. *Preprint*, arXiv:2503.01688.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *Preprint*, arXiv:2406.01574.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.

Chaoyi Wu, Pengcheng Qiu, Jinxin Liu, Hongfei Gu, Na Li, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2025. Towards evaluating and building versatile large language models for medicine. *npj Digital Medicine*, 8(1):58.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 40 others. 2024a. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024b. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*.

Yu Yang, Siddhartha Mishra, Jeffrey N Chiang, and Baharan Mirzasoleiman. 2024c. Smalltolarge (s2l): Scalable data selection for fine-tuning large language models by summarizing training trajectories of small models. *Preprint*, arXiv:2403.07384.

Botao Yu, Frazier N. Baker, Ziqi Chen, Xia Ning, and Huan Sun. 2024. Llasmol: Advancing large language models for chemistry with a large-scale, comprehensive, high-quality instruction tuning dataset. *Preprint*, arXiv:2402.09391.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging LLM-as-a-judge with MT-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023a. Lima: Less is more for alignment. *Preprint*, arXiv:2305.11206.

Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhao Chen. 2023b. Context-faithful prompting for large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14544–14556, Singapore. Association for Computational Linguistics.

# A Prompts

## A.1 Prompts used for complexity estimation

| |
|---|
| *System prompt for a single token response* |
| |
| The following are multiple choice questions about subject. Write down ONLY the NUMBER of the correct answer and nothing else. |
| *System prompt for a single token response with a fallback for unknown answers* |
| |
| The following are multiple choice questions about subject. If you are certain about the answer return the correct option number, otherwise return 0. Write down ONLY the NUMBER and nothing else. |
| *System prompt for a single token response with a fallback for unknown answers (alternative)* |
| |
| The following are multiple choice questions about subject. If you know the answer return the correct option number, otherwise return 0. Write down ONLY the NUMBER and nothing else. |
| *System prompt for a chain-of-thought response* |
| |
| The following are multiple choice questions about subject. Explain your thinking process step-by-step. At the end, write down the number of the correct answer by strictly following this format: [[number of correct answer]]. |
| *System prompt for a chain-of-thought response with a fallback for unknown answers* |
| |
| The following are multiple choice questions about subject. Explain your thinking process step-by-step. At the end, if you are certain about the answer write down the number of the correct answer by strictly following this format: [[number of correct answer]], otherwise return [[0]]. |
| *System prompt for a chain-of-thought response with a fallback for unknown answers (alternative)* |
| |
| The following are multiple choice questions about subject. Explain your thinking process step-by-step. At the end, if you know the answer write down the number of the correct answer by strictly following this format: [[number of correct answers]], otherwise return [[0]]. |
| *User prompt* <br> Question: ... <br> Options: <br> 1. ... <br> 2. ... <br> ... <br> n. ... <br> Choose one of the answers. Write down ONLY the NUMBER of the correct answer and nothing else. |

Table 4: Prompts for complexity estimation

## A.2 General prompts

| |
|---|
| You are an expert in the topic of the question. Please act as an impartial judge and evaluate the complexity of the multiple-choice question with options below. Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, you must not answer the question. You must rate the question complexity by strictly following the criteria: [[Number of reasoning steps]] - how many reasoning steps do you need to answer this question? Valid answers: low, medium, high. Your answer must strictly follow this format: "[[Number of reasoning steps: answer]]". Example 1: "Your explanation... [[Number of reasoning steps: low]]". Example 2: "Your explanation... [[Number of reasoning steps: high]]". Example 3: "Your explanation... [[Number of reasoning steps: medium]]". |

Table 5: Prompt for MASJ reasoning

| |
|---|
| You are an expert in the topic of the question. Please act as an impartial judge and evaluate the complexity of the multiple-choice question with options below. Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, you must not answer the question. You must rate the question complexity by strictly following the scale: "high school and easier", "undergraduate", "graduate", "postgraduate". You must return the complexity by strictly following this format: "[[complexity]]", for example: "Your explanation... Complexity: [[undergraduate]]", which corresponds to the undergraduate level. |

Table 6: Prompt for MASJ education levels

# B Aggregation Methods

Our analysis considers diverse methods to identify the complexity of a questions based on outputs logits with a specific focus on commonly used entropy. They are listed in the list with more details provided below:

1. CoT word-aggregation methods

   - Single Token Answer Entropy
   - CoT Mean
   - CoT Max
   - Difference between CoT Max and Single Token Answer Entropy

2. CoT sequence-aggregation methods

   - Sequence Mean of Words Mean
   - Sequence Max of Words Mean
   - Sequence Mean of Words Max

3. Probability-based methods

   - Mean of Marginal Difference - mean of difference between top-2 probabilities for each token of response

- Top-2 Entropy Difference - difference of top-2 highest entropies for response

4. Hybrid method

- Mix of CoT word-aggregation methods - linear combination of the best perform methods

## B.1 Word-aggregation Methods

This CoT aggregations have the same entropy values as in Section 3.2.3 for each CoT token:

$$h_j = -\sum_{i=1}^{n} p_i \log p_i,$$

where $p_i$ is the probability of a specific token, $n$ is the vocabulary size, and $h_j$ is the entropy of the corresponded token. Aggregating per-token entropies for an answer of length $N$ get:

$$\text{CoT}_{\text{mean}} = \frac{1}{N} \sum_{j=1}^{N} h_j,$$

$$\text{CoT}_{\text{max}} = \max_j h_j.$$

So, Chain-of-Thought maximum and answer entropy difference is:

$$|\text{CoT}_{\text{max}} - h_{\text{answer}}|,$$

where $h_{\text{answer}}$ is the entropy of the answer token.

## B.2 Sequence-aggregation Methods

For $M$ logical claims, which were split by tokens that corresponded to the end of the sequence, we have tokens sets $C_1, C_2, \ldots, C_M$.

$$\text{Seq}_{\text{mean}} = \frac{1}{M} \sum_{j=1}^{M} \left[ \frac{1}{|C_j|} \sum_{i \in C_j} h_i \right],$$

$$\text{Seq}_{\text{mean,max}} = \frac{1}{M} \sum_{j=1}^{M} \left[ \max_{i \in C_j} h_i \right],$$

$$\text{Seq}_{\text{max,mean}} = \max_j \left[ \frac{1}{|C_j|} \sum_{i \in C_j} h_i \right].$$

## B.3 Probability-based Methods

Assume that for each token in response, we have the token probability distribution $p_i$. So, the marginal token difference is $\sigma_i$ and mean marginal difference is mean of all $\sigma_i$ in LLM response.

$$\sigma_i = p_i^{(}1) - p_i^{(}2),$$

$$\overline{\sigma} = \frac{1}{N} \sum_{i=1}^{N} \sigma_i,$$

here $p_i^{(}k)$ is the $k$-th top value in the vector $\mathbf{p}_i$ of the probability distribution of the tokens from the dictionary during the generation of $i$-th token.

We can also consider differences between top two entropies in response $\delta$:

$$\delta = \max_j h_j - \max_{i|i \neq j} h_i.$$

## B.4 Hybrid Method

We can also use a linear combination of 2 best-perform previous methods: $h_{\text{answer}}$ and $\text{CoT}_{\text{max}}$. Also, we tried adding the third element $\text{CoT}_{\text{mean}}$, but it has decreased the ROC-AUC, so we made a decision to remove it.

$$h_{\text{mix}} = (1-\alpha)h_{\text{answer}} + \alpha\text{CoT}_{\text{max}},$$

where $0 \leq \alpha \leq 1$ is the hyperparameter. Empirically we identified that the best value for $\alpha$ is 0.05 for Qwen-3B.

## B.5 Detailed Results

## C Additional experiments

### C.1 MASJ education level and reasoning score

Table 9 shows ROC AUC values for MASJ evaluations of education levels and reasoning scores.

We can see that MASJ reasoning score has a slightly higher ROC AUC of 0.55 on average compared to education levels with ROC AUC of 0.52. There is no significant difference between prompts that allow IDK answers and the ones that do not.

The results indicate that MASJ scores divide the data into complexity groups with moderate quality. On the other hand, results depend on encoding of nominal scores provided by MASJ, and a more comprehensive study could improve this method.

**Technical details.** To calculate ROC AUC we encode MASJ results on a scale from 0 to 1 and prompt the model to answer questions directly, using prompts. For education levels, we take "High school and easier" - 0.2, "Undergraduate" - 0.4, "Graduate" - 0.6, "Postgraduate" - 0.8. For reasoning scores, "Low" - 0.25, "Medium" - 0.5, "High" - 0.75. IDK responses and results with invalid formatting are excluded from the calculations.

### C.2 Feature importances of reasoning model

We evaluated logistic regression weights, that reflect feature importance. Our classifier achieves

| Category | Qwen 3B | Qwen 3B* | Phi4-mini | Phi4-mini* | Phi4 | Phi4* | Mistral 24B | Mistral 24B* |
|---|---|---|---|---|---|---|---|---|
| All | 0.72/0.33 | 0.70/0.33 | 0.72/0.40 | 0.74/0.46 | **0.80**/0.51 | **0.80**/0.58 | 0.75/0.49 | 0.74/0.60 |
| Law | 0.63/0.24 | 0.60/0.21 | 0.64/0.29 | 0.62/0.30 | 0.69/0.47 | 0.69/0.48 | 0.69/0.41 | **0.75**/0.56 |
| Business | 0.67/0.28 | 0.71/0.26 | 0.67/0.31 | 0.64/0.38 | 0.73/0.36 | **0.75**/0.44 | 0.69/0.40 | 0.68/0.43 |
| Psychology | 0.77/0.51 | 0.75/0.51 | **0.84**/0.57 | 0.82/0.59 | **0.84**/0.74 | **0.84**/0.74 | 0.79/0.66 | 0.75/0.68 |
| Chemistry | 0.69/0.23 | 0.62/0.24 | 0.62/0.34 | 0.64/0.41 | 0.70/0.34 | **0.77**/0.45 | 0.68/0.38 | 0.75/0.59 |
| Biology | 0.79/0.59 | 0.79/0.56 | 0.85/0.67 | 0.85/0.73 | **0.90**/0.80 | **0.90**/0.83 | 0.81/0.74 | 0.73/0.80 |
| History | 0.66/0.36 | 0.63/0.35 | 0.68/0.39 | 0.65/0.43 | **0.76**/0.62 | 0.73/0.63 | 0.69/0.54 | 0.64/0.56 |
| Other | 0.70/0.33 | 0.67/0.34 | 0.72/0.39 | 0.74/0.43 | 0.81/0.57 | **0.82**/0.58 | 0.79/0.52 | 0.75/0.59 |
| Physics | 0.65/0.27 | 0.64/0.28 | 0.65/0.32 | 0.66/0.40 | 0.75/0.39 | **0.78**/0.46 | 0.74/0.38 | 0.71/0.63 |
| Computer science | 0.76/0.29 | 0.70/0.32 | 0.73/0.41 | 0.76/0.46 | 0.77/0.55 | **0.80**/0.57 | 0.77/0.51 | 0.74/0.64 |
| Health | 0.69/0.39 | 0.66/0.39 | 0.71/0.43 | 0.71/0.47 | **0.78**/0.64 | 0.77/0.65 | 0.75/0.61 | 0.71/0.63 |
| Economics | 0.77/0.44 | 0.74/0.43 | 0.79/0.55 | 0.80/0.59 | **0.85**/0.68 | 0.83/0.72 | 0.77/0.62 | 0.75/0.66 |
| Math | 0.69/0.24 | 0.67/0.24 | 0.65/0.27 | 0.69/0.31 | 0.73/0.37 | **0.74**/0.43 | 0.69/0.33 | 0.72/0.44 |
| Philosophy | 0.66/0.33 | 0.70/0.31 | 0.71/0.39 | 0.70/0.43 | **0.77**/0.61 | 0.76/0.63 | 0.71/0.53 | 0.70/0.56 |
| Engineering | 0.67/0.34 | 0.66/0.32 | 0.62/0.39 | 0.64/0.45 | 0.74/0.43 | 0.67/0.53 | 0.70/0.46 | **0.77**/0.60 |
| Education level | | | | | | | | |
| High school and easier | 0.73/0.35 | 0.72/0.34 | 0.76/0.38 | 0.75/0.51 | 0.81/0.50 | **0.82**/0.54 | 0.75/0.48 | 0.70/0.52 |
| Undergraduate | 0.73/0.34 | 0.71/0.34 | 0.72/0.42 | 0.77/0.44 | 0.81/0.52 | **0.82**/0.62 | 0.77/0.50 | 0.74/0.64 |
| Graduate | 0.66/0.28 | 0.65/0.26 | 0.64/0.35 | 0.68/0.37 | 0.74/0.50 | 0.73/0.54 | 0.71/0.46 | **0.76**/0.58 |
| Postgraduate | 0.63/0.18 | 0.52/0.20 | 0.64/0.20 | 0.63/0.22 | **0.67**/0.40 | 0.65/0.41 | 0.62/0.35 | 0.63/0.39 |
| MASJ reasoning score | | | | | | | | |
| Low | 0.72/0.42 | 0.71/0.42 | 0.78/0.48 | 0.79/0.51 | 0.82/0.64 | **0.83**/0.65 | 0.79/0.59 | 0.73/0.59 |
| Medium | 0.72/0.32 | 0.70/0.31 | 0.70/0.39 | 0.72/0.44 | **0.79**/0.50 | **0.79**/0.59 | 0.74/0.47 | 0.76/0.63 |
| High | 0.64/0.27 | 0.62/0.27 | 0.59/0.33 | 0.58/0.36 | **0.69**/0.41 | 0.64/0.29 | 0.64/0.39 | 0.62/0.45 |

Table 7: ROC AUC/accuracy for single token response entropy
* Alternative prompt to allow model answer "I do not know"

| Category | Qwen 3B | Qwen 3B* | Phi4-mini | Phi4-mini* |
|---|---|---|---|---|
| All | **0.68**/0.41 | 0.67/0.41 | 0.61/0.43 | 0.58/0.55 |
| Law | **0.60**/0.24 | 0.57/0.23 | 0.55/0.26 | 0.52/0.28 |
| Business | **0.68**/0.45 | 0.67/0.47 | 0.66/0.55 | 0.56/0.65 |
| Psychology | **0.73**/0.51 | 0.70/0.51 | 0.68/0.48 | 0.65/0.63 |
| Chemistry | 0.65/0.41 | **0.68**/0.39 | 0.65/0.43 | 0.63/0.60 |
| Biology | **0.77**/0.56 | 0.68/0.60 | 0.65/0.48 | 0.67/0.71 |
| History | **0.62**/0.36 | 0.61/0.36 | 0.59/0.37 | 0.51/0.39 |
| Other | **0.63**/0.38 | **0.63**/0.36 | 0.60/0.42 | 0.58/0.52 |
| Physics | **0.68**/0.42 | 0.67/0.41 | 0.62/0.39 | 0.61/0.57 |
| Computer science | 0.68/0.37 | **0.73**/0.33 | 0.59/0.41 | 0.58/0.58 |
| Health | **0.63**/0.37 | 0.61/0.40 | 0.62/0.33 | 0.56/0.50 |
| Economics | **0.70**/0.48 | 0.68/0.50 | 0.61/0.47 | 0.65/0.63 |
| Math | **0.73**/0.51 | **0.73**/0.48 | 0.63/0.58 | 0.60/0.67 |
| Philosophy | 0.63/0.33 | 0.62/0.35 | **0.66**/0.37 | 0.59/0.48 |
| Engineering | 0.63/0.31 | **0.64**/0.33 | 0.60/0.37 | 0.55/0.45 |
| Education level | | | | |
| High school and easier | **0.72**/0.56 | 0.70/0.53 | 0.66/0.57 | 0.60/0.73 |
| Undergraduate | **0.67**/0.41 | **0.67**/0.41 | 0.62/0.42 | 0.60/0.55 |
| Graduate | **0.61**/0.27 | 0.60/0.28 | 0.58/0.30 | 0.57/0.38 |
| Postgraduate | 0.63/0.22 | **0.66**/0.15 | 0.41/0.22 | 0.41/0.20 |
| MASJ reasoning score | | | | |
| Low | **0.69**/0.49 | 0.67/0.49 | 0.64/0.46 | 0.61/0.62 |
| Medium | **0.67**/0.41 | 0.66/0.41 | 0.60/0.43 | 0.57/0.55 |
| High | **0.65**/0.26 | 0.60/0.26 | 0.53/0.32 | 0.54/0.36 |

Table 8: ROC AUC/accuracy for single token response entropy after chain-of-thought
* Alternative prompt to allow model answer "I do not know"

| Model | Education level | Reasoning |
|---|---|---|
| Qwen 3B | <u>0.53</u> | 0.55 |
| Qwen 3B* | <u>0.53</u> | 0.55 |
| Phi4-mini | 0.52 | 0.55 |
| Phi4-mini* | 0.52 | 0.54 |
| Phi4 | 0.50 | <u>0.57</u> |
| Phi4* | 0.50 | 0.55 |
| Mistral 24B | 0.50 | 0.56 |
| Mistral 24B* | 0.52 | 0.53 |

Table 9: ROC AUC for MASJ

\* Alternative prompt to allow model answer "I do not know"

| Statistics | Importance |
|---|---|
| Thinking total entropy | <u>1.45</u> |
| Thinking number of tokens | <u>1.08</u> |
| Answer total entropy | 0.25 |
| Answer length | 0.20 |

Table 10: Absolute values of parameter weights

$0.721$, $0.717$ and $0.731$ accuracies by using thinking total entropy, length of the reasoning chain or both features combined correspondingly, thus producing a reasonable model for the evaluation of a probability of the error for question.

Table 10 shows that total entropy and number of tokens of the reasoning chain are the most important parameters influencing the correctness of the model's prediction.

**Technical details.** To avoid excessively long reasoning chains, we set a maximum generation length of 5000 tokens. We also use normalized parameters to remove the mean and scale to unit variance. We take model coefficients of the corresponding parameters as their importance.