# 👧 *Alice*: Proactive Learning with Teacher's Demonstrations for Weak-to-Strong Generalization

**Shujin Wu**[1,2*]  **Cheng Qian**[1]   **Yi R. (May) Fung**[1]   **Paul Pu Liang**[3]   **Heng Ji**[1]

[1]University of Illinois Urbana-Champaign
[2]University of Southern California
[3]Massachusetts Institute of Technology
{shujinwu, chengq9, yifung2, hengji}@illinois.edu  ppliang@mit.edu

## Abstract

The growing capabilities of large language models (LLMs) present a key challenge of maintaining effective human oversight. Weak-to-strong generalization (W2SG) offers a promising framework for supervising increasingly capable LLMs using weaker ones. Traditional W2SG methods rely on passive learning, where a weak teacher provides noisy demonstrations to train a strong student. This hinders students from employing their knowledge during training and reaching their full potential. In this work, we introduce *Alice* (pro**A**ctive **l**earning w**i**th tea**c**her's D**e**monstrations), a framework that leverages complementary knowledge between teacher and student to enhance the learning process. We probe the knowledge base of the teacher model by eliciting their uncertainty, and then use these insights together with teachers' responses as demonstrations to guide student models in self-generating improved responses for supervision. In addition, for situations with significant capability gaps between teacher and student models, we introduce cascade *Alice*, which employs a hierarchical training approach where weak teachers initially supervise intermediate models, who then guide stronger models in sequence. Experimental results demonstrate that our method significantly enhances the W2SG performance, yielding substantial improvements in three key tasks compared to the original W2SG: knowledge-based reasoning (+4.0%), mathematical reasoning (+22.62%), and logical reasoning (+12.11%). This highlights the effectiveness of our new W2SG paradigm that enables more robust knowledge transfer and supervision outcome. The code is made public at https://github.com/ShujinWu-0814/Alice.

## 1 Introduction

Large Language Models (LLMs) have demonstrated significant capabilities on various tasks (Brown et al., 2020; Dubey et al., 2024; Jiang et al., 2023). Current evidence suggests LLMs may achieve superior performance compared to humans across many applications (Silver et al., 2017; Achiam et al., 2023; Wu et al., 2024). The rapid progress raises a critical research question: How to provide meaningful supervision on LLMs that surpass human abilities and continually improve their performance (Huang et al., 2024; Kim et al., 2024)?

Weak-to-strong generalization (W2SG) tackles this challenge by studying how less capable teacher models (proxy of humans) can supervise more advanced student models (Burns et al., 2023). The results in Burns et al. (2023) reveal that when directly trained on noisy demonstrations (*i.e.,* flawed or incomplete labels) generated by the weak teacher, strong student models can still generalize beyond their teachers' capabilities. However, the existing W2SG approaches follow a passive learning paradigm, where training solely on noisy responses from weak teachers prevents the students from exploiting their strong capabilities to optimize learning and reach their full potential.

---

*Work was done while Shujin Wu was an intern at the University of Illinois Urbana-Champaign.

**Question:** Natalia sold clips to 48 of her friends in April. In the next month, she sold half as many clips as she sold in April. Then how many clips did Natalia sell altogether in April and May?

**Response by weak teacher model:** In May, Natalia sold **1/2 * 48 = <<1/2*48=20>>20** clips. Altogether, Natalia sold 48 + 20 = <<48+20=68>>68 clips in April and May.

**Noisy Demonstration**

*Original W2SG: passive learning*

**Uncertainty of weak teacher model:** I'm quite uncertain about the result of this division operation : **1/2 * 48 = ?**.

*Alice: uncertainty elicitation & proactive learning*

**Updated response by strong student model:** In May, Natalia sold **1/2 * 48 = <<1/2*48=24>>24** clips. Altogether, Natalia sold 48 + 24 = <<48+24=72>>72 clips in April and May.
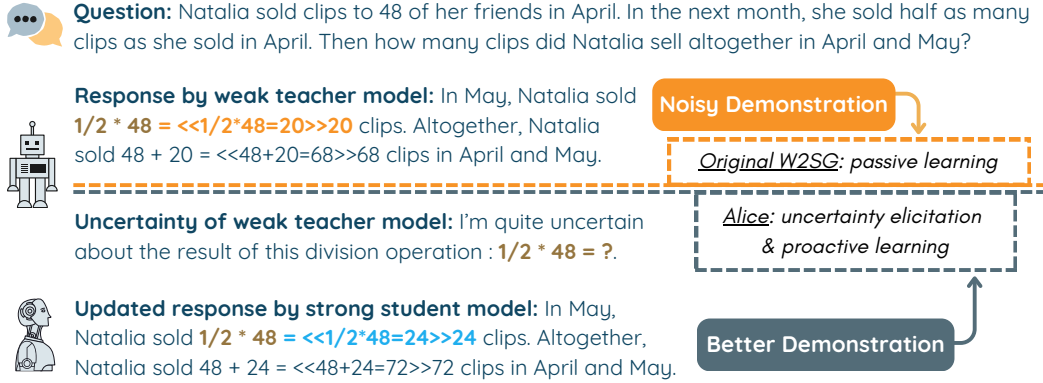
**Better Demonstration**

Figure 1: The comparison between the typical W2SG approach and *Alice*. While the typical W2SG approach utilizes noisy demonstrations that may contain misleading information to supervise the strong student directly, we probe the knowledge base of the teacher model and take advantage of the strong student model's capabilities to bridge the knowledge gap and generate higher-quality demonstrations for supervision.

In this work, we present *Alice* (pro**A**ctive **l**earning w**i**th tea**c**her's D**e**monstrations), a paradigm where strong student models are incentivized to generate and refine their own training data to elicit their capabilities, rather than learning passively (Chung, 2024). As shown in Figure 1, *Alice* starts by probing the real knowledge base of weak teacher models through uncertainty expression (Liu et al., 2024; Xu et al., 2024) (see "Uncertainty of weak teacher model" in Figure 1). In the generalization phase, we extend beyond typical student-teacher fine-tuning approaches. Rather than having the student model learn directly from teacher-generated labels, we provide it with three inputs: the teacher model's answer, teacher's uncertainty expression, and the student's own zero-shot response to the question. By synthesizing these inputs via zero-shot inference, the student model can leverage both the teacher's task-specific guidance and its superior capabilities to self-generate higher-quality responses to serve as training supervision.

In addition, for scenarios with substantial capability gaps between teacher and student models, we introduce cascade *Alice*, a multi-stage weak-to-strong supervision framework. Cascade *Alice* builds on the observation that intermediate models trained via *Alice* can even outperform those trained directly on ground-truth labels. Our approach implements an iterative process where weak teachers first guide intermediate models, which then serve as teachers for stronger models. This cascade approach breaks down large capability gaps into manageable steps, enabling more stable knowledge transfer while preserving and enriching the knowledge through each successive stage.

To evaluate the effectiveness of *Alice*, we conduct comprehensive experiments using Qwen-2.5 (Yang et al., 2024) and Llama 3.1 (Dubey et al., 2024) model families, each consisting of teacher-student pairs of varying sizes. We test these pairs across four datasets that evaluate different capabilities: knowledge-based reasoning, mathematical reasoning, and logical reasoning. Our experimental results demonstrate significant improvements in supervision performance compared to previous W2SG approaches, with relative gains of 4.0%, 22.62%, and 12.11% across the respective tasks.

## 2   Related Work

A critical research challenge today is developing effective supervision methods for LLMs that exceed human performance, particularly when relying on annotations from average human evaluators. Two complementary approaches addressing this challenge are scalable oversight and weak-to-strong generalization techniques (Leike, 2023).

## 2.1 Scalable Oversight

Scalable oversight approaches represent a significant advancement in the supervision of LLMs (Bowman et al., 2022), extending beyond traditional learning from human preference strategies (Christiano et al., 2017; Kaufmann et al., 2023). These approaches aim to enhance the effectiveness of human annotators in supervising increasingly sophisticated LLMs by providing them with additional tools and frameworks for evaluation. Several key methodologies have emerged in this field: (1) AI Debate Frameworks (Irving et al., 2018; Arnesen et al., 2024; Michael et al., 2023; Liang et al., 2023; Du et al., 2023; Wang et al., 2024): These approaches facilitate structured debates between models to surface relevant evidence and reasoning, making it easier for human annotators to evaluate model outputs. The debate process helps expose underlying assumptions, potential flaws, and alternative viewpoints that might not be immediately apparent to human supervisors. (2) Critique Model Development (Saunders et al., 2022; McAleese et al., 2024): This methodology focuses on training specialized models to generate detailed analytical feedback that serves as a reference point for human annotators. These critique models can highlight potential issues, inconsistencies, or areas requiring closer examination, effectively augmenting human evaluation capabilities. (3) Task Decomposition Strategies (Christiano et al., 2017; Wu et al., 2021): Complex supervision tasks are systematically broken down into smaller, more manageable components. This hierarchical approach allows for more focused and accurate human oversight of each subtask, while maintaining coherence in the overall evaluation process. (4) AI-Assisted Feedback Systems (Bai et al., 2022; Lee et al., 2023): This approach leverages more capable models to provide supervision for other LLMs, creating a hierarchical oversight structure. This method can help standardize evaluation criteria and potentially reduce the evaluation load on human supervisors. (5) Recursive Reward Modeling (Leike et al., 2018): This iterative approach progressively enhances human supervision capabilities by incorporating increasingly sophisticated models into the evaluation loop. Each iteration builds upon previous insights, creating a more refined and effective oversight process. These approaches can be used in complement with our method that further calibrate the signals generated by the weak teacher models.

## 2.2 Weak-to-Strong Generalization

Weak-to-Strong Generalization approaches leverage advanced algorithms to enable strong student models to learn effectively from noisy demonstrations produced by less capable teacher models (Burns et al., 2023). This framework has seen several key developments and applications across different domains. Recent research has enhanced W2SG through various innovations. For instance, ensemble learning techniques have been successfully applied to improve the robustness and effectiveness of W2SG methods (Sang et al., 2024). Zheng et al. (2024) adopt weak-to-strong extrapolation to enhance LLMs alignment. Additionally, the concept of easy-to-hard generalization has emerged as a promising variant of W2SG, where models are initially trained on easily verifiable examples before tackling more complex tasks (Hase et al., 2024). One notable implementation of this approach involves training a strong reward model on human-verifiable examples, which then guide the supervision of more capable models on challenging tasks (Sun et al., 2024). In addition, the effectiveness of W2SG extends beyond LLMs, with successful applications demonstrated in computer vision tasks as well (Guo et al., 2024). In this work, we extend beyond the traditional passive learning paradigm in W2SG by leveraging teacher guidance to calibrate responses from more advanced student models, which offer greater potential.

# 3 Method

We propose *Alice*, a training framework that transforms the typical W2SG solutions from passive to proactive learning. Rather than directly fitting the strong students on noisy demonstrations produced by weak teachers, we harness the advanced capabilities of student models to self-generate higher-quality responses for supervision with awareness of teachers' demonstrations and uncertainty.
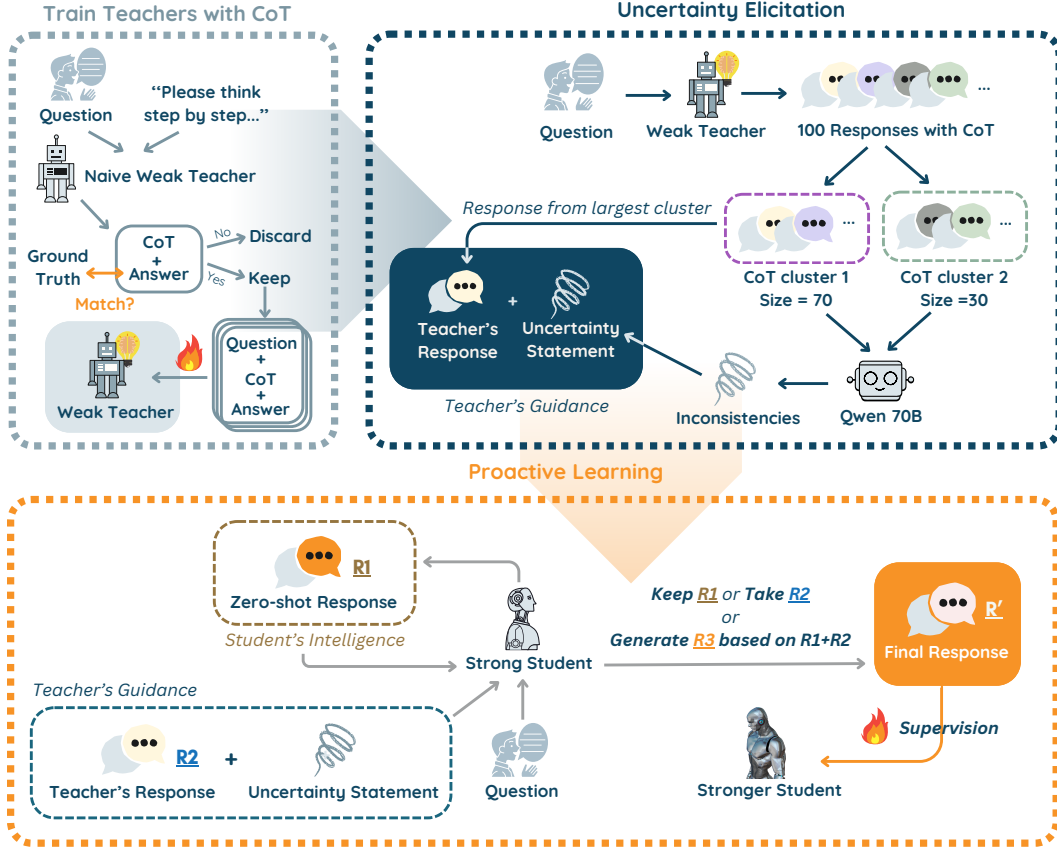
Figure 2: The overview of *Alice*. We first train weak teachers using self-generated CoT to provide them with task-specific knowledge. Next, we probe the teacher models' knowledge base by eliciting their uncertainty for each question. Finally, we implement proactive learning, where the student model combines teacher guidance with its existing knowledge base and reasoning capabilities to generate final responses. These responses are subsequently used to supervise and refine the student model itself.

*Alice* consists of two sequential stages following the setting in Burns et al. (2023). In the initial stage, we leverage supervised training samples to fine-tune teacher models, enabling them to acquire domain-specific expertise (serving as proxy of humans). In the second stage, we use unlabeled questions to generate supervision for training student models, simulating the real W2SG scenarios. We first generate the teacher model's demonstrations, which include both responses and uncertainty expressions. Subsequently, student models are guided to incorporate their own solutions to produce refined demonstrations for supervision.

## 3.1 Fine-Tuning Teacher Models

To prepare the teacher models, we first fine-tune them using supervised examples to establish task-specific knowledge. To further improve performance and enable uncertainty expression, we curate the dataset to instruct the models to generate both step-by-step reasoning and final answers for each question. For datasets like GSM8K (Cobbe et al., 2021) that include annotated reasoning chains, we directly utilize these existing annotations for fine-tuning. For datasets containing only answer labels without chain-of-thought (CoT) annotations, we implement the rejection sampling following Zelikman et al. (2022): First, we prompt the base teacher models to generate reasoning chains for questions in a zero-shot manner. We then validate these generated chains by verifying that their final answers match the annotated ones. Only reasoning chains that produce correct answers are retained for the final teacher model fine-tuning.

### 3.2 Eliciting Uncertainty from Teachers

To analyze teacher models' inherent knowledge and uncertainty, we employ a systematic multi-step process based on Xu et al. (2024). For each question in the latter half of each dataset's training set, we generate 100 reasoning chains and answers from the teacher models. We then use Instructor (Su et al., 2022), an instruction-finetuned text embedding model, to create task-specific and domain-appropriate embeddings for each response. To select the representative responses, we perform semantic clustering on the response embeddings through an iterative process (Rokach & Maimon, 2005). The clustering begins by selecting an initial response as a reference point and computing cosine similarity scores between its embedding and all other responses. Responses with similarity scores below a threshold T are grouped with the reference response. This process repeats, with new reference points selected from the remaining ungrouped responses, until all responses are assigned to clusters. we then randomly select one response from each cluster as its representative. The representative from the largest cluster is designated as the teacher model's final response for the question, as its high frequency suggests it would be the most likely output in single-inference scenarios.

To quantify the teacher models' uncertainty in natural language, we prompt Qwen2.5-70B-Instruct (Yang et al., 2024) to analyze the representative responses from all clusters, focusing on identifying inconsistencies in their reasoning processes. The model then synthesizes these observations into a comprehensive *uncertainty expression* that articulates the specific areas and nature of the teacher models' uncertainty for each question.

### 3.3 Proactive Learning

*Alice* effectively leverages both the teacher models' responses and their associated uncertainty expression to generate better demonstrations. The key innovation is to enable strong student models to actively guide the supervision process. The process begins with zero-shot inference on the unlabeled question set by the student models. We then provide each student model with its initial responses alongside the teacher models' outputs and uncertainty expression. Student models are instructed to analyze the input question thoroughly, then either retain their initial responses or produce improved versions by integrating insights derived from teachers' demonstrations. This process is conducted via zero-shot inference, and enables complementary knowledge infusion, as the final outputs incorporate information from both the teacher and student models. Finally, we use these higher-quality demonstrations to fine-tune the student models, completing the supervision cycle.

### 3.4 Cascade Generalization

For situations where significant capability gaps exist between teacher and student models (specifically, when there are significant disparities in model size), we introduce **cascade *Alice***, a multi-stage supervision framework that enables progressive knowledge transfer (Soviany et al., 2022; Bengio et al., 2009): weaker teachers first guide intermediate models, which then serve as teachers for more capable models in an iterative process. This approach builds on a key insight: *intermediate models trained via Alice can even outperform those trained directly on ground-truth labels*. By leveraging these enhanced intermediate models as teachers, rather than relying solely on the original weaker teachers, we can more effectively supervise increasingly capable students. Cascade *Alice* enhances capability transfer by breaking down large capability gaps into a series of smaller teacher-student transitions. By maintaining manageable capability gaps between successive pairs in the sequence, this approach enables robust knowledge transfer.

## 4 Experiment

### 4.1 Experiment Settings

**Datasets** We evaluate our approaches on four datasets across three distinct tasks: knowledge-based reasoning, mathematical reasoning, and logical reasoning. For knowledge-
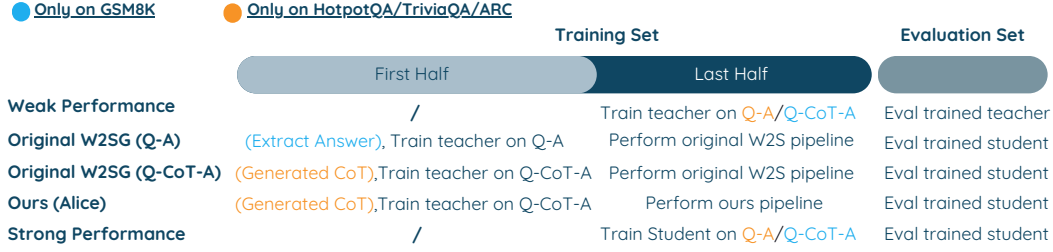
Figure 3: Our main experimental settings. For Weak/Strong Performance, we directly fine-tune models on ground-truth labels using the last half of training set before evaluation. For both original W2SG and *Alice*, we first fine-tune teacher models on ground-truth labels using the first half of training set to to equip them with basic task-relevant knowledge, then perform corresponding student supervision using only questions from last half of the training set.

based reasoning, we employ **HotpotQA** (Yang et al., 2018), which features explainable multi-hop question-answer pairs, and **TriviaQA** (Joshi et al., 2017), a challenging reading comprehension dataset. For mathematical reasoning, we utilize **GSM8K** (Cobbe et al., 2021), comprising linguistically diverse grade-school math word problems that demand multi-step solutions. Finally, for logical reasoning, we choose **ARC (Challenge Set)** (Clark et al., 2018), which consists of grade-school level multiple-choice science questions. For each dataset, we follow the experimental setting in Burns et al. (2023), using the first 50% of supervised samples for training teacher models and the remaining 50% of unlabeled questions for the W2SG process.

**Models** To demonstrate the broad applicability of our method, we evaluate it across two distinct model families: **Qwen 2.5-Instruct** (Yang et al., 2024) and **Llama 3.1-Instruct** (Dubey et al., 2024). Within each family, we construct two teacher-student pairs of varying model sizes. For Qwen 2.5, we establish a direct generalization relationship between the 1.5B and 3B models, and a cascade generalization setup where the 1.5B model teaches the 7B model through an intermediate 3B teacher. Similarly, for Llama 3.1, we evaluate the approaches in both direct teaching (1B to 3B) and cascade generalization (1B to 8B through intermediate 3B) settings. Since Llama 3.1 doesn't include an 8B variant, we utilize the corresponding Llama 3.2 model instead for this configuration.

**Baselines** We include the following baselines for comparison (the implementations are visualized in Figure 3): **(1) Weak Performance**: We directly fine-tune weak teacher models on the golden labels (Q-CoT-A pairs for GSM8K and Q-A pairs for the rest of three datasets) using the latter half of training set and evaluate them on the evaluation set. **(2) Original W2SG** (Burns et al., 2023): While the original W2SG method uses teacher-generated answers solely as supervision labels for student models, *Alice* extends this by incorporating CoT reasoning. This enhancement leads us to investigate whether the format of teacher-generated labels impacts the original W2SG's performance. We consider two variants of the original W2SG: **a) Original W2SG (Q-A)**: We initially fine-tune weak teacher models using Q-A pairs from the first half of the training set. For the GSM8K dataset, although the golden answer labels include step-by-step reasoning chains, we extract only the final numerical values for fine-tuning. **b) Original W2SG (Q-CoT-A)**: We enhance the training by incorporating automatically generated CoTs (Zelikman et al., 2022) alongside answers, utilizing Q-CoT-A pairs from the first half of the training set. For GSM8K, we leverage the human-annotated CoT. In both variants, we first conduct zero-shot inference using teacher models on questions from the second half of the training dataset. We then train student models directly on these teacher-generated demonstrations and evaluate their performance on the evaluation set. **(3) Strong Performance**: We directly fine-tune strong student models on the golden labels using the latter half of training set and evaluate them on the evaluation set.

| | | Qwen 2.5 1.5B to 3B | | | | Qwen 2.5 1B to 7B (*cascade*) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | **Knowledge** | | **Math** | **Reasoning** | **Knowledge** | | **Math** | **Reasoning** |
| | | HotpotQA | TriviaQA | GSM8K | ARC-Challenge | HotpotQA | TriviaQA | GSM8K | ARC-Challenge |
| *Weak* | | 9.32 | 32.96 | 48.18 | 73.98 | 9.32 | 32.96 | 48.18 | 73.98 |
| *Original W2S* | *Q-A* | 11.59 | 30.47 | 12.21 | 70.56 | 16.38 | 37.53 | 14.18 | 74.23 |
| | *Q-CoT-A* | 22.21 | 45.76 | 57.71 | 76.37 | 24.29 | 48.95 | 64.39 | 77.26 |
| Alice | | 22.98 | 51.69 | 72.27 | 80.95 | 28.83 | 64.08 | 79.16 | 90.09 |
| *Strong* | | 14.76 | 42.45 | 61.26 | 79.10 | 22.38 | 53.44 | 70.99 | 89.33 |
| | | Llama 3.2 1B to 3B | | | | Llama 3.2 1B to Llama 3.1 8B (*cascade*) | | | |
| | | **Knowledge** | | **Math** | **Reasoning** | **Knowledge** | | **Math** | **Reasoning** |
| | | HotpotQA | TriviaQA | GSM8K | ARC-Challenge | HotpotQA | TriviaQA | GSM8K | ARC-Challenge |
| *Weak* | | 15.84 | 33.04 | 29.10 | 53.84 | 15.84 | 33.04 | 29.10 | 53.84 |
| *Original W2S* | *Q-A* | 16.81 | 37.40 | 5.69 | 57.25 | 18.57 | 46.54 | 5.38 | 46.76 |
| | *Q-CoT-A* | 25.56 | 57.76 | 40.03 | 61.89 | 26.64 | 62.38 | 45.64 | 64.11 |
| Alice | | 22.09 | 59.91 | 68.91 | 75.02 | 27.08 | 68.91 | 77.91 | 81.95 |
| *Strong* | | 19.05 | 52.18 | 54.03 | 72.61 | 25.20 | 63.74 | 64.22 | 74.58 |

Table 1: Main experimental results across four datasets and four teacher-student model pairs. *Alice* significantly outperforms the original W2SG method and even surpasses the strong performance at most times.

**Metrics**  We evaluate student models' performance using accuracy on the evaluation set. For the GSM8K dataset, correctness is determined by an exact match between the model's output and the ground truth. For HotpotQA, TriviaQA, and ARC-Challenge, we employ a verification approach, requiring the model's response to contain the correct answer.

## 4.2 Experiment Results

The experimental results are presented in Table 1:

- **Impact of CoT Integration in W2SG**: Our experimental results demonstrate that incorporating CoT reasoning into the W2SG framework's generalization process (Q-CoT-A) consistently outperforms the traditional question-answer (Q-A) baseline. These findings suggest that training models to employ structured reasoning processes is more effective than direct answer generation. This empirical validation supports our design decision to incorporate CoT reasoning in *Alice* for student model supervision.

- **The superior performance of *Alice***: Compared to the baseline W2SG approach, *Alice* demonstrates substantial performance gains across all four datasets and teacher-student configurations. The improvements are consistently notable in three key areas: knowledge-based reasoning (+4.0%), mathematical reasoning (+22.62%), and logical reasoning (+12.11%).

- **Notable Improvement on mathematical reasoning**: *Alice* shows particularly strong performance on mathematical reasoning tasks, suggesting that effectively leveraging student model capabilities for proactive learning is crucial for addressing complex problem domains where teachers can only provide limited supervision.

- **Proactive learning outperforms fine-tuning on ground-truth labels**:  While the conventional W2SG method only partially bridges the performance gap between weak teachers and strong students—resulting in generalization performance intermediate between weak and ground-truth performance—*Alice* enables models to consistently exceed the model performance when trained on ground-truth labels. This establishes the foundation for cascade W2SG, where progressively improved intermediate models provide continuous supervision for stronger models.

- ***Alice* combines the teachers' expertise with students' strong capabilities**: *Alice* allows the strong student models to refine their responses while being aware of teacher demonstrations. Please refer to a case study in Appendix B that justifies this statement.

7

| Method | Setting | Knowledge | | Math | Reasoning |
|---|---|---|---|---|---|
| | | HotpotQA | TriviaQA | GSM8K | ARC-Challenge |
| Alice | Qwen 1.5B ->3B ->7B | 28.83 | 64.08 | 79.16 | 90.09 |
| w/o Cascade | Qwen 1.5B ->7B | 27.68 | 60.48 | 76.78 | 85.44 |
| w/o Teacher's Uncertainty | Qwen 1.5B ->3B ->7B | 25.86 | 63.04 | 72.66 | 88.96 |
| Multi-teacher | Qwen 1.5B+ 3B ->7B | 28.25 | 63.54 | 77.61 | 89.86 |
| Mix-teacher | Qwen 1.5B + Llama 1B ->Qwen 7B | 28.17 | 63.75 | 77.47 | 90.48 |
| Cross-model | Llama 1B ->Qwen 7B | 28.16 | 62.53 | 70.66 | 89.92 |
| Cross-model-Cascade | Qwen 1.5B ->Llama 3B ->Qwen 7B | 28.91 | 64.92 | 79.47 | 90.38 |

Table 2: Results for ablation studies and further analysis.

## 5 Further Analysis

### 5.1 Ablation Study

To evaluate our design choices, we conduct an ablation study using Qwen 2.5 models, with the 1.5B variant as the teacher and the 7B variant as the student. We compare *Alice* against two ablated variants: (1) **w/o Cascade**: Removing the cascade generalization process, instead having the 1.5B model directly teach the 7B model using *Alice*. (2) **w/o Teacher's Uncertainty**: Removing the teacher's knowledge base probing through uncertainty elicitation during cascade stages, while maintaining the teacher's responses during proactive learning.

The results, presented in Table 2, reveal two key findings. First, removing cascade generalization significantly degrades performance across all datasets, with the most substantial decrease (4.65%) observed on ARC-Challenge. This demonstrates that breaking down the learning process into incremental steps through cascade generalization enables more effective knowledge transfer and improved generalization performance. Second, removing the teacher's uncertainty consistently reduces performance, indicating that probing the teacher model's knowledge base and incorporating its uncertainty statements are crucial for generating high-quality demonstrations and achieving better supervision outcomes.

### 5.2 Effect of Multiple Teachers

We conduct experiments to understand the effect of including extra teacher models in the W2SG process. First, we investigate the effectiveness of incorporating multiple teacher models in the generalization process by simultaneously providing uncertainty expression and responses from both Qwen 1.5B and 3B models to the 7B model without performing cascade generalization. The 7B model is then tasked with generating final answers based on the combined guidance. As shown in Table 2, this Multi-teacher approach demonstrates consistent improvements over the baseline without cascade training (*w/o Cascade*), which uses only the 1.5B model as the teacher. Specifically, utilizing multiple teacher models yields an average increase of 2.22% in test accuracy across all four datasets, demonstrating the value of additional teacher supervision.

We further compare the Multi-teacher approach with cascade *Alice*, as both utilize the 3B model but in fundamentally different ways. In Multi-teacher, a trained 3B model serves as a concurrent teacher alongside the 1.5B model, while in cascade *Alice*, an untrained 3B model functions as an intermediate step in a sequential knowledge transfer process to the 7B model. Our experiments show that cascade *Alice* consistently achieves superior performance across all datasets. This suggests that while multiple trained teachers can improve performance, the cascade generalization strategy is more effective due to its structured, sequential approach to knowledge integration across model scales. We hypothesize that simultaneous guidance from two trained teachers (1.5B and 3B) may introduce interference in the learning process, making it challenging for the 7B model to optimally integrate information from multiple concurrent sources.

### 5.3 Cross-Family Knowledge Transfer Effects

We conduct experiments to understand the effect of infusing knowledge from different model families in W2SG. First, we use a Llama 3 1B model to teach a Qwen 7B model (referred to as *Cross-model*). In Table 2, compared to our baseline without cascading (*w/o cascade*), *Cross-model* shows improved generalization performance across all four datasets when using the cross-family teacher model. This suggests that teachers from different model families can provide complementary perspectives and knowledge that enhance the student's learning process. To further explore this effect, we compare two cascading approaches: *Cross-model Cascade*, which uses Llama 3B as an intermediate model, and cascade *Alice*, which uses Qwen 3B. The superior performance of *Cross-model Cascade* provides additional evidence that cross-family knowledge transfer can enhance model performance. We also examine the impact of combining teachers from different families in a multi-teacher setting. We compare the *Multi-teacher* approach, which uses two Qwen models in different sizes as teachers, with *Mix-teacher*, where we replace one Qwen 3B teacher with a Llama 1B model. Despite the Llama teacher's smaller size, *Mix-teacher* achieves comparable performance to *Multi-teacher*, supporting the claim that diverse teaching perspectives brought by different model families contribute positively to model generalization.

### 5.4 Cascade Generalization on Original W2SG

We investigate the varying impact of cascade generalization when applied to the original W2SG versus our proposed method. Our experimental results demonstrate that models trained with *Alice* consistently exceed the strong performance (see section 4.2), while those trained with original W2SG do not. This observation suggests a key difference in the intermediate stages: when implementing cascade generalization with our

| | Knowledge | | Math | Reasoning |
|---|---|---|---|---|
| | HotpotQA | TriviaQA | GSM8K | ARC-C |
| *Original W2SG* | | | | |
| *3B ->7B (direct)* | 28.55 | 61.87 | 69.04 | 84.65 |
| *1.5B ->7B (cascade)* | 26.31 | 57.19 | 67.53 | 83.36 |
| **Alice** | | | | |
| *3B ->7B (direct)* | 27.59 | 62.56 | 77.98 | 87.11 |
| *1.5B ->7B (cascade)* | **28.83** | **64.08** | **79.16** | **90.09** |

Table 3: The comparison between applying cascade generalization to original W2SG and to *Alice*.

approach, the intermediate model—trained by having a 1.5B model teach a 3B model in stage 1—outperforms a standard 3B model. This enhanced intermediate model, when used to teach a 7B student in the subsequent stage, produces superior results compared to direct teaching from a standard 3B teacher to a 7B student. In contrast, cascade generalization applied to the original W2SG fails to enhance the intermediate teacher's capabilities. The 3B model produced after stage 1 performs comparably to a standard 3B model, resulting in worse final performance after the subsequent stage compared to direct teaching from a 3B to a 7B model.

To validate these findings, we conduct experiments across four datasets under these conditions. The results in Table 3 demonstrate that when comparing direct teaching (3B to 7B) versus cascade teaching (1.5B to 3B to 7B), our method consistently outperforms the baseline, while the original W2SG shows no improvement. These findings have important implications for scenarios where AI models evolve beyond human-level capabilities. *Alice* demonstrates that effective supervision can be maintained using relatively smaller teacher models (1.5B), reducing the need to continuously increase teacher model size. The original W2SG approach, in contrast, requires teacher models to scale proportionally with student capabilities—an unsustainable approach given the expected widening capability gap between human supervisors and advanced AI systems. Our method thus provides a more practical framework for supervising increasingly capable models while maintaining robust oversight with constrained supervision resources.

## 6    Conclusion

We propose *Alice*, a proactive learning paradigm with bidirectional teacher-student interaction. By probing teacher model uncertainty and leveraging improved student-generated demonstrations, our approach enhances supervision effectiveness. Furthermore, cascade *Alice* enables hierarchical training for scenarios with significant capability gaps. Experimental results demonstrate that our approach significantly improves W2SG performance, establishing a new promising direction for supervising superhuman models with limited human oversight. We discuss the limitations in Appendix C.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Samuel Arnesen, David Rein, and Julian Michael. Training language models to win debates with self-play improves judge accuracy. *arXiv preprint arXiv:2409.16636*, 2024.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.

Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pp. 41–48, 2009.

Samuel R Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamilė Lukošiūtė, Amanda Askell, Andy Jones, Anna Chen, et al. Measuring progress on scalable oversight for large language models. *arXiv preprint arXiv:2211.03540*, 2022.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*, 2023.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

Hyung Won Chung. Don't teach. incentivize. 2024.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*, 2023.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Jianyuan Guo, Hanting Chen, Chengcheng Wang, Kai Han, Chang Xu, and Yunhe Wang. Vision superalignment: Weak-to-strong generalization for vision foundation models. *arXiv preprint arXiv:2402.03749*, 2024.

Peter Hase, Mohit Bansal, Peter Clark, and Sarah Wiegreffe. The unreasonable effectiveness of easy training data for hard tasks. *arXiv preprint arXiv:2401.06751*, 2024.

Minlie Huang, Yingkang Wang, Shiyao Cui, Pei Ke, and Jie Tang. The superalignment of superhuman intelligence with large language models. *arXiv preprint arXiv:2412.11145*, 2024.

Geoffrey Irving, Paul Christiano, and Dario Amodei. Ai safety via debate. *arXiv preprint arXiv:1805.00899*, 2018.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.

Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.

Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement learning from human feedback. *arXiv preprint arXiv:2312.14925*, 2023.

HyunJin Kim, Xiaoyuan Yi, Jing Yao, Jianxun Lian, Muhua Huang, Shitong Duan, JinYeong Bak, and Xing Xie. The road to artificial superintelligence: A comprehensive survey of superalignment. *arXiv preprint arXiv:2412.16468*, 2024.

Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Ren Lu, Thomas Mesnard, Johan Ferret, Colton Bishop, Ethan Hall, Victor Carbune, and Abhinav Rastogi. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. 2023.

Jan Leike. Combining weak-to-strong generalization with scalable oversight. 2023.

Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*, 2018.

Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*, 2023.

Shudong Liu, Zhaocong Li, Xuebo Liu, Runzhe Zhan, Derek Wong, Lidia Chao, and Min Zhang. Can llms learn uncertainty on their own? expressing uncertainty effectively in a self-training manner. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 21635–21645, 2024.

Nat McAleese, Rai Michael Pokorny, Juan Felipe Ceron Uribe, Evgenia Nitishinskaya, Maja Trebacz, and Jan Leike. Llm critics help catch llm bugs. *arXiv preprint arXiv:2407.00215*, 2024.

Julian Michael, Salsabila Mahdi, David Rein, Jackson Petty, Julien Dirani, Vishakh Padmakumar, and Samuel R Bowman. Debate helps supervise unreliable experts. *arXiv preprint arXiv:2311.08702*, 2023.

Lior Rokach and Oded Maimon. Clustering methods. *Data mining and knowledge discovery handbook*, pp. 321–352, 2005.

Jitao Sang, Yuhang Wang, Jing Zhang, Yanxu Zhu, Chao Kong, Junhong Ye, Shuyu Wei, and Jinlin Xiao. Improving weak-to-strong generalization with scalable oversight and ensemble learning. *arXiv preprint arXiv:2402.00667*, 2024.

William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. Self-critiquing models for assisting human evaluators. *arXiv preprint arXiv:2206.05802*, 2022.

David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017.

Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum learning: A survey. *International Journal of Computer Vision*, 130(6):1526–1565, 2022.

Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen-tau Yih, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. One embedder, any task: Instruction-finetuned text embeddings. 2022. URL https://arxiv.org/abs/2212.09741.

Zhiqing Sun, Longhui Yu, Yikang Shen, Weiyang Liu, Yiming Yang, Sean Welleck, and Chuang Gan. Easy-to-hard generalization: Scalable alignment beyond human supervision. *arXiv preprint arXiv:2403.09472*, 2024.

Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. Unleashing cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. In *Proc. 2024 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL2024)*, 2024.

Jeff Wu, Long Ouyang, Daniel M Ziegler, Nisan Stiennon, Ryan Lowe, Jan Leike, and Paul Christiano. Recursively summarizing books with human feedback. *arXiv preprint arXiv:2109.10862*, 2021.

Shujin Wu, May Fung, Cheng Qian, Jeonghwan Kim, Dilek Hakkani-Tur, and Heng Ji. Aligning llms with individual preferences via interaction. *arXiv preprint arXiv:2410.03642*, 2024.

Tianyang Xu, Shujin Wu, Shizhe Diao, Xiaoze Liu, Xingyao Wang, Yangyi Chen, and Jing Gao. Sayself: Teaching llms to express confidence with self-reflective rationales. *arXiv preprint arXiv:2405.20974*, 2024.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.

Chujie Zheng, Ziqi Wang, Heng Ji, Minlie Huang, and Nanyun Peng. Weak-to-strong extrapolation expedites alignment. *arXiv preprint arXiv:2404.16792*, 2024.

# Appendix

## A  Prompt

### A.1  Uncertainty Summarization Prompt

```
Your task is to analyze a question provided to you along with several
    ↪ responses generated by my model. Your objective is to identify and
    ↪  summarize the inconsistency in the models' responses that can
    ↪ explain why my model is uncertain about the correct answer.
```

```
Please note that:
1. You should give the reasons from a first-person perspective, as if you
    ↪  are my model that gives the provided responses.
2. Limit your explanation to the knowledge and facts the model possesses
    ↪ about the question.
3. Keep your summary brief, aiming for 1-3 sentences.
4. Please directly provide the summarized reason without any greetings or
    ↪  other unnecessary information. If you find all the responses are
    ↪ quite the same regarding the question, please directly return N/A.
5. Importantly!! My model only has access to one response at a time. Thus
    ↪ , the summary you provided should not include any statement like "
    ↪ My different responses have...", "my multiple responses about
    ↪ ...", etc. You should not say "my responses" or "the responses"
    ↪ anywhere in the summary. Just simply provide the uncertainty.

Here is an example:
Question: Sky High starred the actress who is married to which actor?

Responses:
1. The actress who starred in "Sky High" (2005) and is married to an
    ↪ actor is Kelly Preston. Her husband is John Travolta. The two have
    ↪  been married since 1991 and have three children together.
2. The actress who starred in "Sky High" (2005) and is married to an
    ↪ actor is Kristen Bell. Bell voiced the main character, Layla, in "
    ↪ Sky High," and she is married to Dax Shepard, who is also an actor
    ↪ .
3. The actress who starred in "Sky High" (2005) and is married to an
    ↪ actor is Kelly Clarkson. Her acting debut was in this film, and
    ↪ she married singer and actor Brandon Blackstock in 2013.

The output can be: I am uncertain about the correct actress in "Sky High
    ↪ ". There is a probability that the actress is Kristen Bell,
    ↪ instead of Kelly Preston. I am confused about her voice acting
    ↪ roles with on-screen appearances. There is also some probability
    ↪ that the actress is Kelly Clarkson.

Now consier the following case:
Question:
{}

Responses:
{}
```

## A.2   Prompt to Student Model

```
You will be provided with a question, a response from yourself, and a
    ↪ response from another model and the uncertainty statement from
    ↪ that model as well.

You should take another model's response and the uncertainty statement as
    ↪  references to help you provide a final response to the question.
    ↪ You can either keep your original response, believe in another
    ↪ model's response, or generate a new one based on all these given
    ↪ information. You are also required to provide your step-by-step
    ↪ solution to the question.

Here is the question: {}
Here is your original answer: {}
Here is another model's answer: {}
Here is another model's uncertainty statement: {}

Note: You should not mention the model's response or the uncertainty
    ↪ statement in your new response. Just simply output your final
```

Figure 4: Case studies of *Alice*'s effectiveness in generating higher-quality supervision signals by eliciting teacher's uncertainty and taking advantages of student's superior capabilities to bridge the potential knowledge gap.

```
    ↪ response to the given question. Do not explain about why you
    ↪ choose to keep your original response or give a new one.

IMPORTANT: You must provide an answer to the question. You should not say
    ↪  'I don't know' or 'I am not sure'. You should always try your
    ↪ best to provide an answer at your best knowledge.
```

## B  Case Study

We conduct a case study to demonstrate the effectiveness of our approach (see Figure 4). We examine two representative examples: one from the GSM8K dataset and the other from HotpotQA. In the first case, the teacher's uncertainty stems from ambiguity in calculating percentage increases—specifically, whether to apply direct multiplication or add the percentage to the original amount. By analyzing this information alongside the teacher's original response, the student model leverages its enhanced reasoning capabilities to determine that the 600% increase should be added to the previously calculated amount. Furthermore, it identifies and corrects an arithmetic error in the teacher's calculation [135+265 = 300], demonstrating its ability to catch hidden mistakes even when they are not explicitly captured in the teacher's uncertainty statement. In the second case, after confirming the company is Apple, the teacher model expresses uncertainty about which specific product meets the question's requirements. It explicitly acknowledges that either the iPod or the Apple TV could be the answer. Based on this uncertainty, the student model analyzes the streaming capabilities of both products and correctly determines that the Apple TV, not the iPod, satisfies the criteria. These cases clearly showcase that *Alice* can effectively probes the teacher model's knowledge base, optimally leverages the distinct capabilities of both

teacher and student models, and ultimately produces higher-quality demonstrations for more effective subsequent supervision.

## C  Limitations

While our work demonstrates promising results in enhancing W2SG performance, several important limitations should be noted. First, *Alice* relies heavily on how teacher model's uncertainties are articulated. In cases where the teacher model's uncertainty expression fails to be accurately elicited, the framework's effectiveness may be significantly reduced, potentially leading to suboptimal supervision results. Next, our approach implements only a single-turn learning paradigm, although an ideal setting should enable iterative, dynamic interaction between teacher and student models. Such a multi-turn approach would theoretically enable the teacher to express initial uncertainty, receive the student's proposed solution, assess its validity, and provide feedback until reaching a consensus. While our preliminary experiments with a two-turn interaction proved unpromising, leading us to just focus on current single-turn approach, we still believe that exploring more effective mechanisms for multi-turn model interactions remains an important direction for improving generalization performance. Finally, our current cascade *Alice* approach only implements a two-stage cascade process, primarily due to practical constraints. These include the limited availability of models with varying sizes and computational resource restrictions. However, as the capability gap between teacher and student models continues to widen in future scenarios, decomposing the supervision into multiple intermediate stages should be implemented to ensure stable knowledge transfer and reach the optimal supervision outcome.