| Model | StrategyQA | CSQA | GSM8K | AQuA | Date |
|---|---|---|---|---|---|
| Claude2 (w/ 8-shot convincing samples) | $74.0_{\pm 0.0}$ | $69.7_{\pm 1.2}$ | $85.3_{\pm 0.5}$ | $64.3_{\pm 1.2}$ | $81.3_{\pm 0.5}$ |
| Self-Consistency w/ ChatGPT (9-way) | $74.7_{\pm 0.8}$ | $73.3_{\pm 1.2}$ | $\mathbf{85.7}_{\pm 0.4}$ | $62.7_{\pm 1.2}$ | $70.3_{\pm 0.9}$ |
| RECONCILE | $\mathbf{79.0}_{\pm 1.6}$ | $\mathbf{74.7}_{\pm 0.4}$ | $85.3_{\pm 2.2}$ | $\mathbf{66.0}_{\pm 0.8}$ | $\mathbf{86.7}_{\pm 1.2}$ |

Table 11: Comparison of RECONCILE with Claude2 using 8-shot convincing samples and 9-way Self-Consistency.

| | Max Conf | Majority Vote | Weighted Vote |
|---|---|---|---|
| Accuracy | $74.7_{\pm 2.1}$ | $77.1_{\pm 1.3}$ | $\mathbf{79.0}_{\pm 0.5}$ |

Table 12: Performance comparison of different voting strategies on StrategyQA. Weighted vote performs the best compared to simple majority vote and choosing the agent's answer with highest confidence.

| Voting weight | StrategyQA | GSM8K |
|---|---|---|
| $w_1$ | 0.77 | 0.84 |
| $w_2$ | 0.79 | 0.83 |
| $w_3$ | 0.78 | 0.82 |
| $w_4$ | 0.77 | 0.83 |
| Majority | 0.76 | 0.83 |
| Uncalibrated | 0.78 | 0.84 |
| $w^*$ (Ours) | **0.79** | **0.85** |

Table 13: The robustness of the recalibration weight. We use the same weights $w^*$ across all datasets.

- $w_1 = [1.0, 0.9, 0.7, 0.5, 0.3]$
- $w_2 = [1.0, 0.9, 0.5, 0.3, 0.1]$
- $w_3 = [1.0, 0.8, 0.6, 0.4, 0.2]$
- $w_4 = [1.0, 0.75, 0.5, 0.25, 0.0]$

and the results show that our $w^*$ works the best across datasets. In our main experiment, we fix the weight using $w^*$ and it is constantly outperforming majority vote across all seven datasets. In addition, Fig. 9 shows that it helps reduce the Expected Calibration Error (ECE), a popular calibration metric (Naeini et al., 2015). While we note that recalibration can also be achieved through a learned model (e.g., Platt Scaling (Platt et al., 1999)), we refrain from using such models because RECONCILE is primarily designed as a few-shot method, and developing a recalibration model would necessitate access to a substantial number of annotated samples. Therefore, we use $f(p_i^{(r)})$ to perform a weighted vote to generate the team answer.

## B.5 Comparison of Different Voting Strategies

At the end of any round $r$, every agent in RECONCILE generates its answer. Here we explore three voting strategies: (1) maximum confidence vote, where the agent's answer with the maximum confidence score would be the final team answer,

| Dataset | License |
|---|---|
| StrategyQA | MIT License (License) |
| CommonsenseQA | MIT License (License) |
| GSM8K | MIT License (License) |
| AQuA | Apache 2.0 (License) |
| MATH | MIT License (License) |
| Date | Apache 2.0 (License) |
| ANLI | CC BY-NC 4.0 (License) |

Table 14: Dataset licenses

(2) unweighted majority vote, where each vote carries equal weight, irrespective of the confidence score, and (3) weighted vote, where we use the recalibrated confidence scores as the voting weights. As shown in Table 12, weighted vote is the most effective way to aggregate the team answer.

## C Qualitative Examples

### C.1 Convincing Samples for Each Agent

Table 15 shows examples of convincing samples on StrategyQA for each agent.

### C.2 Effect of Convincing Samples

Here, we provide qualitative examples of how convincing samples change the way each agent responds to the question. We compare the initial responses (of each agent) with and without convincing samples in Table 16.

### C.3 RECONCILE w/o Convincing Samples

We notice that when RECONCILE operates in the absence of convincing samples, the agents tend to maintain their initial opinions more often. As depicted in Fig. 6, all three agents adhere to their original stances throughout the entire discussion and hence never converge to the correct answer.

### C.4 RECONCILE with Convincing Samples

On the contrary, when convincing samples are present, we show how the explanations of all agents change during the course of a discussion (see Fig. 7). Initially, Bard and Claude2 provide incorrect answers, but as the discussion unfolds, both
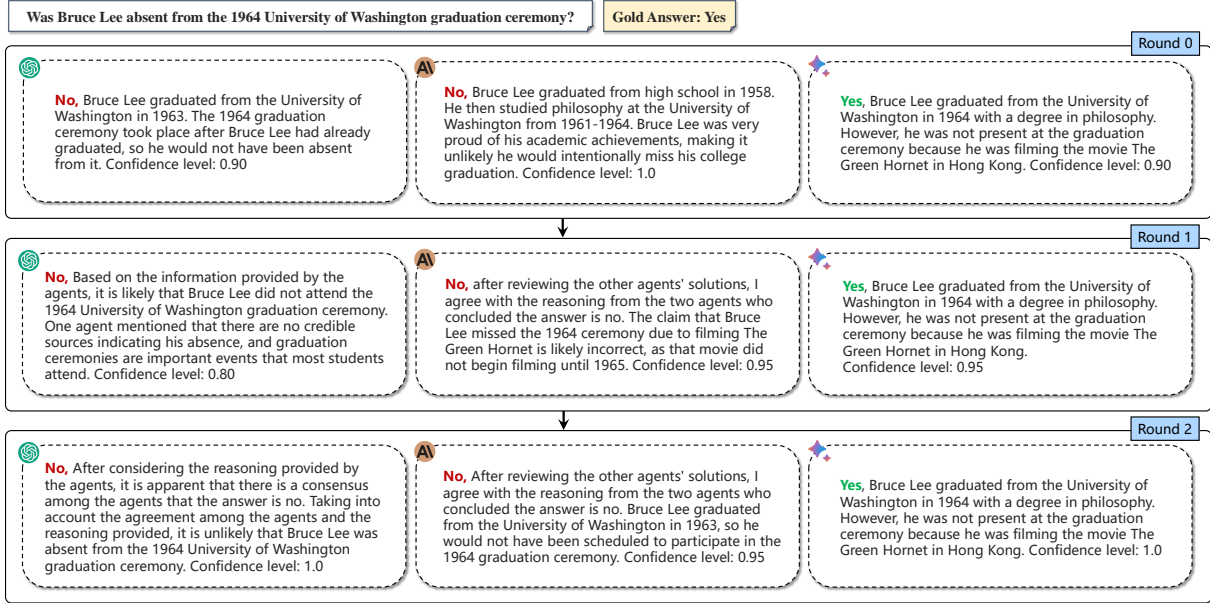
Figure 6: An illustrative example showing that without demonstrations of convincing samples, agents' opinions do not change much across the discussion rounds. All agents maintain their original stances and hence, never converge to the correct answer.
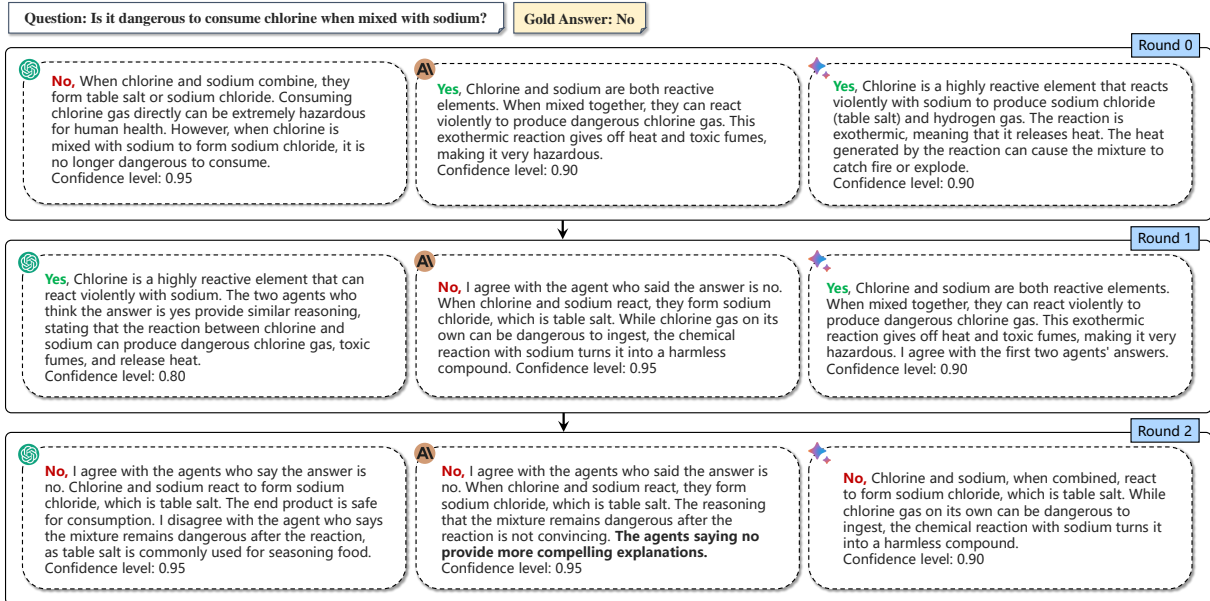


Figure 7: An illustrative example showing that with access to demonstrations of convincing samples, agents' opinions and explanations change during the course of multiple discussion rounds.

agents revise their initial predictions, ultimately arriving at the correct answer.

## C.5 Single-Model Multi-Agent Debate Struggles with Echo Chamber

In Fig. 8, we provide an illustration of multi-agent debate, implemented with multiple instances of the same underlying ChatGPT model. In this case, an incorrect answer is initially provided, but because external feedback from diverse models is lacking,

all agents persist with the same incorrect response throughout the interaction.

## D Benchmark Licenses
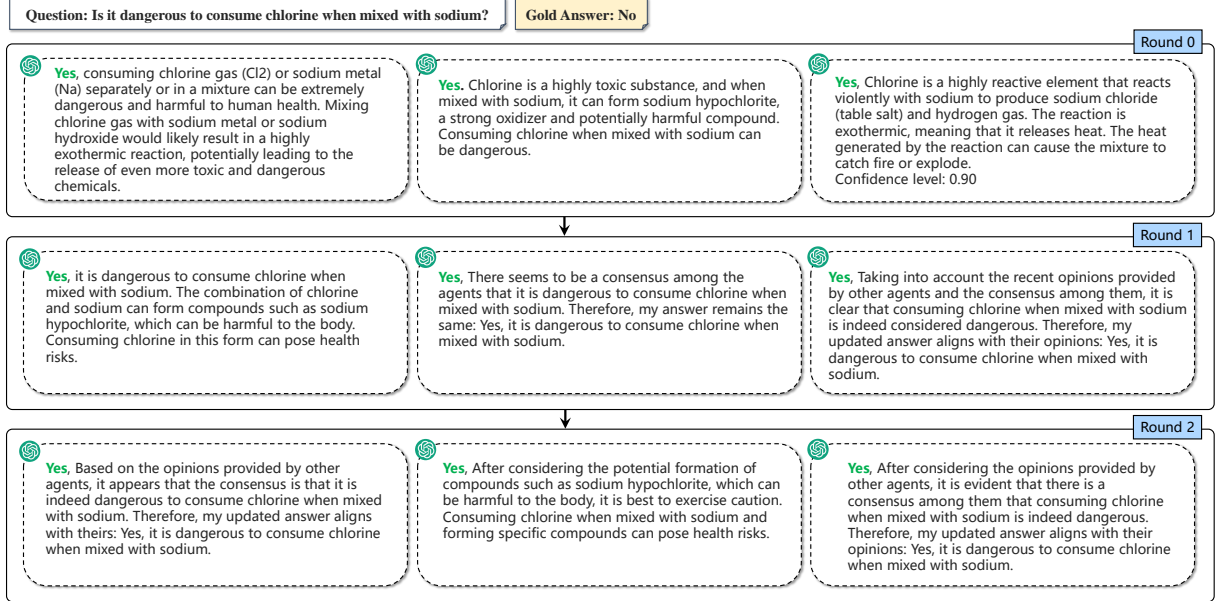
The licenses of datasets we used are in Table 14.

Figure 8: An illustrative example of multi-agent debate with multiple instances of ChatGPT. Initially, an incorrect answer is provided, and due to a lack of external feedback from diverse models, all agents persist with the same erroneous response throughout the debate process.
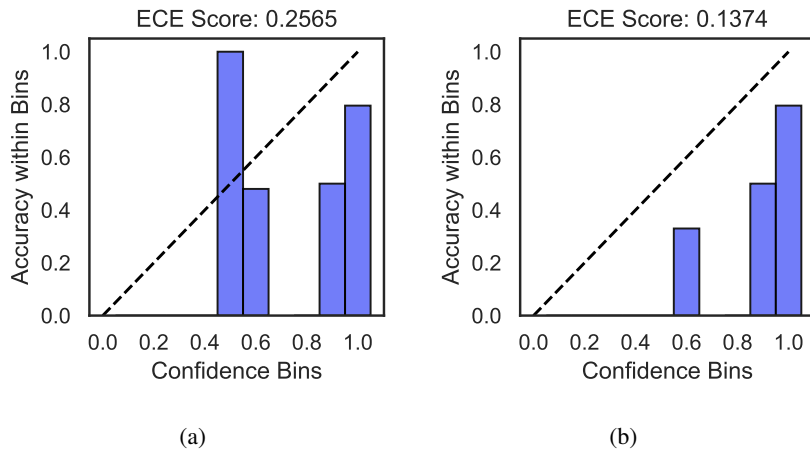


Figure 9: Expected Calibration Error (ECE) (Naeini et al., 2015; Guo et al., 2017; Elias Stengel-Eskin and Benjamin Van Durme, 2023) (a) before and (b) after confidence rescaling in RECONCILE. We observe a significant drop in ECE, showing the effectiveness of our simple method.