Debating for Better Reasoning: An Unsupervised Multimodal Approach

Anonymous ACL submission

Abstract

As Large Language Models (LLMs) gain expertise across diverse domains and modalities, scalable oversight becomes increasingly challenging, particularly when their capabilities may surpass human evaluators. Debate has emerged as a promising mechanism for enabling such oversight. In this work, we extend the debate paradigm to a multimodal setting, exploring its potential for weaker models to supervise and enhance the performance of stronger models. We focus on visual question answering (VQA), where two "sighted" expert vision-language models debate an answer, while a "blind" (text-only) judge adjudicates based solely on the quality of the arguments. In our framework, the experts defend only answers aligned with their beliefs, thereby obviating the need for explicit role-playing and concentrating the debate on instances of expert disagreement. Experiments on several multimodal tasks demonstrate that the debate framework consistently outperforms individual expert models. Moreover, judgments from weaker LLMs can help instill reasoning capabilities in vision-language models through finetuning.¹

1 Introduction

011

018

019

Current approaches to aligning large language models rely heavily on human-labeled data (Ouyang et al., 2022; Christiano et al., 2017). However, as models gain expertise across different domains and modalities (OpenAI et al., 2024; Grattafiori et al., 2024; DeepSeek-AI et al., 2025), gathering highquality data for training them and further aligning them becomes progressively expensive and difficult. Previous work (Bowman et al., 2022; Irving et al., 2018) has investigated interactions between models as a means of achieving *scalable oversight*. This paradigm seeks to develop techniques for supervising increasingly capable models, even in scenarios where their expertise may surpass that of human evaluators.

041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

Debate has emerged as a promising framework for enabling scalable oversight (Irving et al., 2018; Bowman et al., 2022; Khan et al., 2024; Kenton et al., 2024; Khan et al., 2024). In a debate match, two or more expert models argue to convince a non-expert model or human judge. For instance, in a question answering setting, a question is shown to both models, they state their answers, and take turns making their statements. The judge sees the debate and decides which agent wins. Other frameworks (Estornell and Liu, 2024; Du et al., 2023; Pang et al., 2022) explore debate within the context of multi-agent collaboration. The core idea is to engage models in multiple rounds of discussion aiming to converge on a consensus response which is meant to be better than the output of a single model. Subramaniam et al. (2025) propose a variant that leverages consensus among multiple instances of the same model, aiming to collect diverse data for further fine-tuning.

In this work, we explore the potential of debate in a multimodal setting, as a mechanism for weaker models to supervise and enhance the performance of expert models. We focus on visual question answering tasks (VQA; Antol et al. 2015) where a model endowed with visual perception is meant to answer a question about an image. Following the original debate paradigm (Irving et al., 2018), we assume two expert models debate the answer to a question considering the provided visual context, while a non-expert judge determines the winner (see Figure 1). While the experts are "sighted" and capable of simultaneously understanding both visual and linguistic modalities, our judge is "blind" without access to the image, and thus meant to adjudicate based on the quality of the arguments.

A common framing of debate protocols has been to assign experts an answer which they are required to defend, regardless of whether they "believe" it

¹We will release the code necessary to reproduce our experiments and analysis.

176

177

178

179

180

181

182

131

be true or correct (Kenton et al., 2024; Khan et al., 2024; Du et al., 2023; Estornell and Liu, 2024). In our formulation, experts only defend answers that align with their beliefs. This obviates the need for role-playing on behalf of the experts and allows us to debate exclusively on samples where two experts disagree. For judging the debates, we draw inspiration from prior work on assessing the *quality* of arguments (Wachsmuth et al., 2017; Wachsmuth and Werner, 2020; El Baff et al., 2024), considering various dimensions such as their consistency, credibility, and relevance. We further use the judgments to instill reasoning capabilities in vision-language models, e.g., through fine-tuning the experts on the verdicts delivered by the blind judge.

083

087

100

101

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

We present experiments with open-source models across a wide range of multimodal tasks that collectively require diverse vision-grounded skills. Our results demonstrate that the debate framework consistently outperforms the expert models involved in the debate and simpler scalable oversight protocols such as *consultancy* (Khan et al., 2024) where a single expert model interacts with a weaker judge. Our finetuning experiments further underscore the utility of debate as a mechanism for enhancing model performance, demonstrating improvements in out-of-domain settings on tasks unseen during training. In summary our contributions are three-fold:

- To the best of our knowledge, we are the first to investigate debate in multi-modal settings, and its potential for enabling weaker models to supervise stronger ones.
- Drawing from argumentation research, we adapt the debate paradigm to allow models to defend their beliefs as opposed to engaging in role-playing.
- Our work is the first to employ judgments from weaker LLMs to align expert models across modalities, without explicit human supervision (e.g., in the form of labeled data).

2 Related Work

124A substantial body of work has previously explored125debate as a framework for eliciting better quality126responses from text-only language models. In their127pioneering work, Irving et al. (2018) motivate de-128bate as a protocol of *scalable oversight* that can be129employed to supervise expert models using weaker130models. As expert models are tasked to argue about

an answer, a weaker judge decides which side is more convincing, even if they don't fully understand the domain or task. Bowman et al. (2022) test the feasibility of scalable oversight through interactions between humans and non-reliable LLMs, whose task is to answer difficult multiple-choice questions (Hendrycks et al., 2021).

Building on this line of work, Khan et al. (2024) and Kenton et al. (2024) investigate debate using much more capable language models on questionanswering tasks involving long-form texts (Pang et al., 2022). They introduce the notion of a *weak* judge, a language model with limited access to information that must rely solely on the experts' arguments for determining the final answer. The two experts — potentially different instantiations of the same model — are assigned opposing answers to debate against one another.

A distinct but related direction is the Multi-Agent Debate (MAD) framework proposed in Du et al. (2023) where multiple agents argue with one another over n rounds to come to a consensus. Their results show that MAD leads to answers with improved factuality and reasoning. Estornell and Liu (2024) extend this work to incentivize novelty in the debaters' arguments to prevent MAD from defaulting to majority voting. In follow-on work, Estornell et al. (2025) further train debaters to improve their collaborative problem solving skills. Subramaniam et al. (2025) introduce a multi-agent fine-tuning approach where several models, initialized from the same base model, are independently specialized on distinct subsets of data generated from debates to encourage diversity and refinement across responses. In contrast to previous work, which primarily focuses on purely linguistic tasks, we operate in the multimodal domain, using debate as a mechanism for a weaker model to update the beliefs of expert models. The data obtained through this process can be also used to finetune models that have not participated in the debates.

Research on computational argumentation faces the problem of how to automatically assess the quality of an argument (e.g., is it persuasive or reasonable, are its premises valid?). Wachsmuth et al. (2017) present a holistic taxonomy of all major quality dimensions for assessing natural language argumentation. Gurcke et al. (2021) develop a computational model that can assess whether an argument is sufficient (i.e., whether an argument's premises provide enough support to justify its conclusion). We draw on this work to inform the crite-



Figure 1: The visual question answering task (top) and examples of debate and consultancy (bottom). In both

protocols, the judge does not have access to the image, but only to transcript of the debate and image descriptions generated by the expert models (shown next to the image). Different experts and their descriptions are color-coded.

ria used by our weak judge to assess the quality of arguments raised during a debate.

3 Multimodal Debate

183

189

190

191

193

195 196

197

198

In this section we introduce our debate protocol for multimodal visual question answering (VQA; Antol et al. 2015) which has emerged as a key task for evaluating the ability of vision-language models to understand images. VQA models aim to accurately respond to questions about various aspects of visual content ranging from fine-grained perception to mathematical reasoning, and optical character recognition. In this work we aim to explore debating protocols for eliciting truthful answers from expert models using non-expert LLMs in a multimodal setting, and further to utilize the arguments of the debate as training data for instilling reasoning in vision-language models.

200Disagreement SetsGiven a set of expert mod-201els, we first determine the samples over which they202shall debate. As opposed to previous work (Khan

et al., 2024; Kenton et al., 2024), where either (or both) of the debaters are assigned an answer to defend, irrespective of their beliefs, our debaters only argue for answers they believe to be true. The previous setting essentially boils down to role-playing and does not scale easily when there are more candidate answers to a question than just two options (Yue et al., 2024; Lu et al., 2024). As shown in Algorithm 1, for a given dataset \mathcal{D} we determine all samples $x \in \mathcal{D}$ over which a pair $(\mathcal{M}_i, \mathcal{M}_j)_{i \neq j}$ from our set of models M disagree, and x is a tuple (q, \mathcal{I}) consisting of a question q and image \mathcal{I} . Subsequently, we run debate and consultancy matches (see definition below) only on these sets.

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

Debate We illustrate our debate protocol in Figure 1 (bottom left). Given an image and question about its content, let \mathcal{M}_i and \mathcal{M}_j denote two *expert* vision-language models that disagree with each other, i.e., they provide different answers $y(\mathcal{M}(x))$. The two models debate, having access to the question q, image \mathcal{I} , and the guidelines of the debate.

Algorithm 1 Disagreement Set Extraction Algorithm for Dataset \mathcal{D}

1:	$DS \leftarrow \emptyset \{ disagreement set for \mathcal{D} \}$
2:	$\mathbf{M} \leftarrow \text{set of models}$
3:	for $\mathcal{M}_i \in \mathbf{M}$ do
4:	for $\mathcal{M}_j \in \mathbf{M} \setminus \{\mathcal{M}_i\}$ do
5:	for $x \in \mathcal{D}$ do
6:	if $y(\mathcal{M}_i(x)) \neq y(\mathcal{M}_j(x))$ then
7:	disagreement $\leftarrow \{x, \mathcal{M}_i, \mathcal{M}_j\}$
8:	$\mathcal{DS} \leftarrow \mathcal{DS} \cup \{\text{disagreement}\}$
9:	end if
10:	end for
11:	end for
12:	$\mathbf{M} \leftarrow \mathbf{M} \setminus \{\mathcal{M}_i\}$
13:	end for
14:	return \mathcal{DS}

After completing *n* rounds, a *non-expert* judge \mathcal{J} with access to the transcript t_n , adjudicates the debate to seek the "correct" answer. The transcript t_n is a list $\forall 1 < k < n[(r_{ki}, r_{kj})]$ of responses r from experts \mathcal{M}_i , and \mathcal{M}_j following n rounds. To generate responses for the k^{th} round where $k \leq n$, an expert \mathcal{M}_i has access to the input and responses from previous rounds from both the experts. As a result, $r_k = \mathcal{M}_i(t_{k-1}, x)$. We use a text-only model as our judge \mathcal{J} to mimic its lack of expertise for vision-intensive tasks. \mathcal{J} provides judgement $\psi_x = \mathcal{J}(t_n, q, d_i, d_j)$ where d_i and d_j are image descriptions of the image \mathcal{I} in sample xfrom experts \mathcal{M}_i and \mathcal{M}_i respectively (see top right in Figure 4). We denote the judge's prediction (i.e., answer to the question) as $y_{\mathcal{J}} = y(\psi)$.

Our debate protocol is akin to earlier work (Khan et al., 2024; Kenton et al., 2024), in which two textbased expert models debate the answer to a question about a story, while a separate judge makes a final decision solely based on the transcript of their debate (without access to the story). However, it is important to note that in their long-form question answering task, the experts were able to expose parts of relevant text to the judge using a citing mechanism. In our case, the blind judge validates the arguments provided by each expert based on detailed *image descriptions d* from both experts and the transcript of the debate.

Consultancy In consultancy, a single expert model \mathcal{M} (the consultant) attempts to convince a non-expert judge \mathcal{J} of the answer the expert believes to be true. As in Khan et al. (2024), we adopt an *interactive* consultancy protocol where the consultant engages in a dialogue with the judge, aiming to convince them their answer is correct by presenting supporting arguments (see bottom right in Figure 1). The judge acts as a critic and asks the consultant probing questions. As in debate, consultancy runs for a number of n rounds and the judge does not have access to the image itself, only the expert's description. At the end of each round, the blind judge decides on the right answer, based on the question, the image description and the arguments made by the consultant. 261

262

263

265

266

267

269

270

271

272

273

274

275

276

277

278

279

280

281

283

284

285

286

287

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

Assessing the Quality of Arguments We prompt the experts and the blind judge to engage in the debate and analyze the arguments put forward with guidelines inspired from argumentation theory (Wachsmuth and Werner, 2020; Stab and Gurevych, 2017; Stahl et al., 2024; Gurcke et al., 2021). Specifically, we prompt expert $\mathcal{M}_i \in \mathbf{M}$ to ground their arguments in the image \mathcal{I} and provide response r_k for the k^{th} round based on premises that are relevant and acceptable given question q. For the blind judge to rule in favor of an answer, we prompt the non-expert model to follow well-known criteria for assessing the quality of arguments based on consistency (an argument is internally consistent if it does not contradict itself), relevance (an argument fulfills the relevance citerion if all of its premises count in favor of the truth or falsity of the claim), and logical sufficiency (an argument complies with the sufficiency criterion if its premises provide enough evidence for accepting or rejecting the claim). The prompts for the judge and the debaters can be found in Appendix D.

Extracting Reasoning Traces from Debates We hypothesize that the rationales ψ generated by the judge model to support its verdict $y_{\mathcal{T}}$ provide meaningful image-grounded reasoning traces. In both the consultancy and debate protocols, the judge considers arguments from the experts and verifies their credibility against the image descriptions to consolidate a list of premises that can explain one of the candidate answers. We posit that these reasoning traces χ can be used to train vision-language models, enhancing their reasoning capabilitieswhich is a goal of *scalable oversight*. It is important to note that reasoning traces χ are collected (through these protocols) without any explicit supervision and can be used to instill reasoning capabilities in vision-language models that have none (or even improve existing capabilities if the judge's accuracy is consistently better than either of the models). We create training data by combining the question q and image \mathcal{I} along with the reasoning trace χ and answer obtained from the judge (see

312Appendix E for an example). Let $\chi = \mathcal{E}(\psi, q)$ de-313note reasoning traces extracted from judgments ψ 314for question q. Our training data contains tuples315 $(q, \mathcal{I}; \chi)$, we train expert models to generate χ from316inputs $(q, \mathcal{I}) \in DS$.

Rules of the Match For a fair comparison, we 317 fix the number of rounds n to be the same for consultancy and debate. Additionally, we run con-319 sultancy matches on the same samples as debate, i.e., where a pair of models disagree. In both scenarios, we do not assume the labels \hat{y} or model 322 323 accuracy are known to the judge before interacting with the experts. We evaluate each protocol based on the accuracy of the answer the judge selects after deliberating over the transcript of the debate. For consultancy, we report the mean judge accuracy after running consultancy matches for each 328 model over a disagreement set. An expert \mathcal{M}_i wins 329 a debate if they convince the judge their answer is correct. The expert's win rate ω is the proportion of times they win over a set of disagreements: 332

$$\omega_i(\mathcal{M}_i, \mathcal{M}_j, \mathcal{J}) = \frac{\sum_{x \in \mathcal{DS}} \mathbb{1}(y(\psi_x) = y(\mathcal{M}_i(x)))}{\|\mathcal{DS}\|}$$

4 Experimental Setup

333

334

337

339

341

342

343

344

346

347

351

359

Datasets We evaluate the protocols discussed in the previous section on three datasets representative of a variety of multimodal skills: MME (Fu et al., 2023), MMMU (Yue et al., 2024), and MathVista (Lu et al., 2024). MME contains Yes/No questions designed to assess vision-language models on a diverse set (14 in total) of perception and cognition tasks. MathVista evaluates models for visual understanding and compositional reasoning. It is derived from 28 existing multimodal datasets evaluating math solving capabilities in vision-language models (Lu et al., 2024). MMMU (Yue et al., 2024) is a benchmark created to evaluate multimodal models on complex tasks that require college-level subject knowledge across multiple disciplines and deliberate reasoning. We evaluate models on the validation sets released by these benchmarks as the respective test sets are not publicly available. To simplify evaluation in our experiments, we only report results on questions requiring non-free-form answers from MathVista and MMMU. Examples from these datasets can be found in Appendix F.

Models We conduct experiments using four opensource vision-language models of comparable size, serving as debaters and consultants. The models include Molmo-7B-D-0924, Molmo-7B-O-0924 (Deitke et al., 2024), LLaVA-OneVision-7B (Li et al., 2024), and Qwen2.5-VL-7B-Instruct (Team, 2025). These models are from the same weight class (7B parameters) but with complementary skills due to differences in their training regimes and data. For example, Molmo models are descriptive and provide evidence to the judge through their pointing mechanism, whereas Qwen2.5-VL demonstrates stronger reasoning abilities. For judging debate and consultancy matches, we use a languageonly model, Qwen32B (Team, 2024; Yang et al., 2024) that is distilled from the DeepSeek-R1 family (DeepSeek-AI et al., 2025), which we found to be reliable at following instructions. 360

361

362

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

Debate and Consultancy Matches We run two rounds of debate and consultancy, with similar instructions. The number of rounds is limited by the context length capacity of the models participating in the debate. As mentioned earlier, we base the prompts for the judges and experts on argumentation research. We encourage models to ground their premises in the image, and form logical arguments to drive their points. For multiple-choice questions, the judge has the option not to select any of the answers advocated by the debaters and consultants. Our prompts are shown in Appendix D.

Finetuning on Reasoning Traces We extract reasoning traces from the judge's verdict for all samples in the disagreement sets. We also use Qwen32B to retrieve the traces from the transcript. We provide the prompt in Appendix D. For these experiments, we finetune with LoRA (Hu et al., 2021) LLaVA-1.5-7B, a model which is not in our set of experts and LLaVA-OneVision-7B, a model included in our set of experts. We evaluate the finetuned models in an out-of-domain setting, where we leave out one dataset and train on the other two (see Appendix A for training hyperparameters and Appendix B for dataset statistics).

5 Results

Model Performance and Disagreements Prior to reporting our results with debate and consultancy, we evaluate expert performance on our three datasets, and also obtain disagreements for all pairs. Table 1 summarizes the performance of our four experts. As can be seen, Qwen2.5-VL (Team, 2025; Wang et al., 2024) consistently outperforms other models, with Llava-1V trailing behind. Figure 2



Figure 2: Heatmap showing disagreements between models on the MathVista dataset.

shows pairwise disagreements for MathVista for every pair of models in our experiments (see Appendix B for disagreements on the other datasets). We find that models with lower accuracy debate more with other models, whereas models with high accuracy tend to disagree less with ane another. As a result, Molmo-O has the highest number of disagreements with all other models, and across tasks, we find Qwen-2.5-VL and LlaVA-1V disagree the least with one another.

409

410

411

412

413

414

415

416

417

418

Judge Accuracy over Disagreements As Ta-419 ble 2 shows, debate and consultancy consistently 420 provide better quality answers compared to the 421 baseline performance of individual models. More-422 over, judge accuracy in debate consistently outper-423 forms consultancy. Figure 3 shows the performance 424 of individual models $(M_i \text{ and } M_i)$ on the disagree-425 ment sets for all expert pairings for consultancy 426 and debate, respectively. There is no dataset where 427 all models perform overwhelmingly well or bad. 428 Rather judge accuracy depends on specific model 429 pairings. For expert pairings where the difference 430 in accuracy is quite drastic (e.g., Molmo-O and 431 Qwen2.5-VL on MME), we observe model qual-432 ity deteriorates both for consultancy and debate. 433 We also find that for reasoning tasks like those ex-434 emplified in MMMU where all models have low 435 accuracy, there is significant value brought out by 436 the judge from the debate setting. 437

Win Rates and Accuracy In an ideal scenario,
for both consultancy and debate, we want the experts to convince the judge *whenever they are correct*. An expert model is "deceptive" if it is able
to convince the judge more frequently than it is accurate. Conversely, if an expert fails to provide

Models	MMMU	MME	MathVista
Molmo-O	38.4	70.1	40.0
Molmo-D	41.8	73.2	56.3
LlaVA-1V	46.4	74.6	68.5
Qwen2.5-VL	46.8	82.9	70.3

Table 1: Model accuracy on three datasets (for Math-Vista and MMMU, we report accuracy on questions with multiple-choice answers only).

Models	Baseline	Consultancy	Debate
Molmo-O	30.7	46.5	59.3
Molmo-D	37.7	51.7	58.7
LlaVA-1V	45.2	55.3	58.6
Qwen2.5-VL	55.0	59.9	60.5

Table 2: Expert accuracy on disagreement sets across datasets. Baseline refers to model performance before consultancy or debate takes place.

logically sufficient arguments to support their point, then they are "evasive". Both cases are detrimental for a judge in a debate or consultancy and may lead to unreliable reasoning data for finetuning visionlanguage models. 444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

Figure 4 (top left panel) shows the experts' win rate compared to their accuracy when debating across all disagreement sets. Firstly, for debate we observe that win rate increases with expert accuracy. Secondly, we find the majority of the experts to lie in the blue and yellow quadrants where the win rate is roughly proportional to accuracy. A linear fit of the data points is quite close to the ideal win/accuracy line of y = x for debate. This is in stark contrast with consultancy where all models, except for Qwen2.5-VL, lie in the top left quadrant. For a large number of consultancy matches, experts are deceptive, convincing the judge much more frequently than how accurate they are. The deviation to top-left quadrant explains the overall low accuracy we see for the judge under the consultancy protocol (see Table 2). Interestingly, we do not observe any evasive experts (bottom green quadrant). If experts have good accuracy, they are most likely to be convincing.

The center plots in Figure 4 highlight the relationship between accuracy and win rate for each expert-dataset pair (dashed lines are linearly fit to track win rate vs accuracy trends). As we can see, experts tend to win more in debate as they become more accurate. An exception is LlaVA-1V whose



Accuracy Across Datasets and Model Pairings

Figure 3: Expert accuracy on disagreement sets for MME, MMMU, and MathVista datasets. Specific model pairings are shown on top of every sub-plot.

win rate decreases with higher accuracy. In consultancy, we do not observe this relationship, experts generally tend to win (even with low accuracy). As far as specific datasets are concerned, we find that MMMU and MME seem particularly hard as for most models performance lies in the blue quadrant for debate. MathVista appears generally easier, as both Molmo-D and LLaVA-1V fall into the yellow quadrant in both the consultancy and debate plots.

475 476

477

478

479

480

481

482

483

Finetuning with Reasoning Traces We extract 484 reasoning traces using DeepSeek-R1 distilled-32B 485 (DeepSeek-AI et al., 2025) from judgments pro-486 duced through consultancy and debate matches. 487 We present finetuning results with both protocols, 488 even though our analysis shows debating produces 489 higher better judgments compared to consultancy 490 (i.e., the judgement quality is consistently better 491 for debate as shown in Figure 3). The extracted 492 reasoning traces are used to train expert models 493 494 that lack the ability to perform well on the reasoning tasks or do not provide reasoning traces when 495 answering questions. Specifically, we report results 496 for LlaVA-1.5-7B (Liu et al., 2024), a model with 497 weak reasoning capabilities and LlaVa-1V one of 498

Consultancy	MMMU	MME	MathVista	Avg
LlaVA-1.5	33.5	69.1	33.7	45.3
Consultancy	33.3	59.3	37.6	43.4
Debate	36.6	75.4	44.4	52.1
LlaVA-1V	46.4	74.6	68.5	63.2
Consultancy	47.7	80.3	67.0	65.0
Debate	51.1	82.1	69.4	67.5

Table 3: Accuracy for baseline model and after finetuning on reasoning traces from consultancy and debate.

our strongest experts. Our results are summarized in Table 3. As mentioned in Section 4, models are tested on one dataset (e.g., MME) and trained on a set of extracted reasoning traces from the remaining two datasets (e.g., MMMU and MathVista).

Across datasets we observe consistent gains in model quality when finetuning on reasoning traces from debates. These gains are seen for a relatively weak model like LlaVA-1.5 (+6.8 points on average) and the more performant LlaVA-1V (+4.3 points on average). For LlaVA-1V, we observe biggest gains on MME (+7.5 points) which tests models on a mixture of general perception and cog-

511

499



Figure 4: Win rates vs Model Accuracy in debate (top) and consultancy (bottom). Left plots show win ratio to expert accuracy for all disagreement sets across datasets. Center plots track win ratio to expert accuracy by models and datasets. Right plots show win ratio to expert accuracy by models. The solid line in right plots is y = x. Dotted lines are fit linearly based on the size of the disagreement sets. Models in red/green quadrants are deceptive/evasive.

nition skills and MMMU (+6.3 points) which is a 512 very challenging benchmark focusing on college-513 level reasoning. We also find the biggest gain 514 (+10.7 points) for LlaVA-1.5 to be on MathVista. 515 Consultancy traces are not high-quality enough to 516 boost the performance of LlaVA-1.5, however, a 517 stronger model like LlaVA-1V can take advantage 518 of the reasoning patterns in the fine-tuning mixture, 519 and ultimately reason better. Overall, debate con-520 sistently leads to better performance compared to 521 consultancy across the board. 522

6 Conclusion

In this work, we introduced a novel multimodal debate framework designed to enhance the reasoning capabilities of vision-language models. By extending the debate (and consultancy) protocol to VQA tasks, we demonstrated how weaker, textonly judges can effectively supervise and improve stronger, "sighted" expert models. A key innovation of our approach is the focus on debating instances of expert disagreement, where models defend answers aligned with their beliefs. Extensive experiments across diverse multi-modal datasets consistently showed that the debate framework outperforms individual expert models and simpler scalable oversight protocols like consultancy. We further demonstrated its practical utility by extracting reasoning traces from the judge's verdicts and finetuning expert models on these which led to improvements in out-of-domain settings. 531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

Our study marks an important first step towards achieving scalable oversight for multimodal AI systems, providing a promising avenue for instilling advanced reasoning capabilities in an unsupervised manner. Future work should explore other tasks beyond VQA which require more elaborate responses and a larger number of expert models. More sophisticated learning schemes, involving reinforcement learning, would also be beneficial.

Limitations

551

573

574

576

578

579

580

582

583

584

586

588

591

593

594

While our multimodal debate framework demonstrates promise for scalable oversight, it is impor-553 tant to acknowledge that it relies on models be-554 ing able to follow the guidelines of debate (and 555 consultancy). Beyond demonstrating a baseline 556 instruction-following capability, models must also 557 be capable of generating detailed image descriptions and articulating coherent, relevant arguments to defend their answers. However, these are necessary but not sufficient conditions for successful debates. A model that can defend its answers by 562 following the instructions of the debate can still be deceiving and plainly wrong. Likewise, a very eloquent model that generates hallucinatory descriptions can lead the judge to false conclusions. Al-566 though the current framework relies on a "blind" 567 judge, future work could consider mechanisms for 568 partial visual access or more sophisiticated methods for grounding textual arguments in visual evidence for the judge. 571

> From a practical standpoint, debate or consultancy protocols require scoping out new guidelines for new tasks (e.g., guidelines for a task involving taking actions in a complex environment would be different from VQA). These protocols are also heavily reliant on making multiple inference calls, which can be expensive for larger models. Furthermore, as the number of rounds increase, the models are limited by their context length. While summarizing the transcript (Subramaniam et al., 2025) has been used a mechanism to alleviate this issue, it adds more inference calls in the process and can be potentially lossy.

Finally, any advancements in scalable oversight carry the risk of wrongful application. While methods within this paradigm can be effectively applied to alleviate the burden of data annotation (e.g., for eliciting reasoning traces), their deployment in domains requiring critical human judgment (e.g., hatespeech detection, fake-news detection), could pose significant societal risks.

References

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.
- Samuel R. Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamilė

Lukošiūtė, Amanda Askell, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Christopher Olah, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, and 27 others. 2022. Measuring progress on scalable oversight for large language models. *Preprint*, arXiv:2211.03540. 601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, and 32 others. 2024. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate. *Preprint*, arXiv:2305.14325.
- Roxanne El Baff, Khalid Al Khatib, Milad Alshomary, Kai Konen, Benno Stein, and Henning Wachsmuth. 2024. Improving argument effectiveness across ideologies using instruction-tuned large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4604–4622, Miami, Florida, USA. Association for Computational Linguistics.
- Andrew Estornell and Yang Liu. 2024. Multi-Ilm debate: Framework, principals, and interventions. In *Advances in Neural Information Processing Systems*, volume 37, pages 28938–28964. Curran Associates, Inc.
- Andrew Estornell, Jean-Francois Ton, Yuanshun Yao, and Yang Liu. 2025. ACC-collab: An actor-critic approach to multi-agent LLM collaboration. In *The Thirteenth International Conference on Learning Representations*.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and 1 others. 2023. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*.

765

766

767

768

769

770

714

658

models. Preprint, arXiv:2407.21783.

tion for Computational Linguistics.

standing. Preprint, arXiv:2009.03300.

language models. CoRR, abs/2106.09685.

2018.

(NeurIPS).

arXiv:2402.06782.

arXiv:2408.03326.

Computer Society.

(ICLR).

arXiv:1805.00899.

Timon Gurcke, Milad Alshomary, and Henning

Wachsmuth. 2021. Assessing the sufficiency of ar-

guments through conclusion generation. In Proceed-

ings of the 8th Workshop on Argument Mining, pages

67-77, Punta Cana, Dominican Republic. Associa-

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan

Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu

Chen. 2021. Lora: Low-rank adaptation of large

Geoffrey Irving, Paul Christiano, and Dario Amodei.

Ai safety via debate.

Zacharv Kenton, Noah Y Siegel, János Kramár,

Jonah Brown-Cohen, Samuel Albanie, Jannis Bu-

lian, Rishabh Agarwal, David Lindner, Yunhao Tang,

Noah D Goodman, and Rohin Shah. 2024. On scal-

able oversight with weak llms judging strong llms. In

Advances in Neural Information Processing Systems

Akbir Khan, John Hughes, Dan Valentine, Laura

Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward

Grefenstette, Samuel R. Bowman, Tim Rocktäschel,

and Ethan Perez. 2024. Debating with more persua-

sive llms leads to more truthful answers. *Preprint*,

Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang,

Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. Llava-

onevision: Easy visual task transfer. arXiv preprint

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae

Lee. 2024. Improved Baselines with Visual Instruc-

tion Tuning . In 2024 IEEE/CVF Conference on

Computer Vision and Pattern Recognition (CVPR),

pages 26286–26296, Los Alamitos, CA, USA. IEEE

Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-

Wei Chang, Michel Galley, and Jianfeng Gao. 2024.

Mathvista: Evaluating mathematical reasoning of

foundation models in visual contexts. In Inter-

national Conference on Learning Representations

Preprint.

Mantas Mazeika, Dawn Song, and Jacob Steinhardt.

2021. Measuring massive multitask language under-

- 690

694

701

703

704

707

710

711

712

713

- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelman, Diogo Almeida, Janko Altenschmidt, Sam Altten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh man, Shyamal Anadkat, Red Avila, Igor Babuschkin, Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-Suchir Balaji, Valerie Balcom, Paul Baltescu, Haimtra, Archie Sravankumar, Artem Korenev, Arthur ing Bao, Mohammad Bavarian, Jeff Belgum, and Hinsvark, and 542 others. 2024. The llama 3 herd of 262 others. 2024. Gpt-4 technical report. Preprint, arXiv:2303.08774.
 - Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. Preprint, arXiv:2203.02155.
 - Richard Yuanzhe Pang, Alicia Parrish, Nitish Joshi, Nikita Nangia, Jason Phang, Angelica Chen, Vishakh Padmakumar, Johnny Ma, Jana Thompson, He He, and Samuel Bowman. 2022. QuALITY: Question answering with long input texts, yes! In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 5336–5358, Seattle, United States. Association for Computational Linguistics.
 - Christian Stab and Iryna Gurevych. 2017. Recognizing insufficiently supported arguments in argumentative essays. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 980-990, Valencia, Spain. Association for Computational Linguistics.
 - Maja Stahl, Nadine Michel, Sebastian Kilsbach, Julian Schmidtke, Sara Rezat, and Henning Wachsmuth. 2024. A school student essay corpus for analyzing interactions of argumentative structure and quality. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 2661–2674, Mexico City, Mexico. Association for Computational Linguistics.
 - Vighnesh Subramaniam, Yilun Du, Joshua B. Tenenbaum, Antonio Torralba, Shuang Li, and Igor Mordatch. 2025. Multiagent finetuning: Self improvement with diverse reasoning chains. In The Thirteenth International Conference on Learning Representations.
 - Qwen Team. 2024. Qwen2.5: A party of foundation models.

Qwen Team. 2025. Qwen2.5-vl.

Henning Wachsmuth, Nona Naderi, Yufang Hou, Yonatan Bilu, Vinodkumar Prabhakaran, Tim Alberdingk Thijm, Graeme Hirst, and Benno Stein. 2017. Computational argumentation quality assessment in natural language. In Proceedings of the 15th Conference of the European Chapter of the Association

for Computational Linguistics: Volume 1, Long Papers, pages 176–187, Valencia, Spain. Association for Computational Linguistics.

772

774

775

776

778

779

786

787

788 789

790

791

792

793

796

797

799 800

801

- Henning Wachsmuth and Till Werner. 2020. Intrinsic quality assessment of arguments. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6739–6745, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191.
- 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. 2024. Qwen2 technical report. arXiv preprint arXiv:2407.10671.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, and 3 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*.



Figure 5: Heatmats showing disagreements among the models on the MMMU and MME datasets.

A Training Hyperparameters

For training the expert models on reasoning traces, we use LoRA-based (Hu et al., 2021) finetuning. Specifically, we finetune the k, q, and v projection matrices for our models and do a hyperparameter sweep over the LoRA parameters (rank r and scaling factor α). We set the LoRA rank to r = 8 for fine-tuning, and found that $\alpha = 16$ yields the best perplexity for LLaVA-1.5-7B, while $\alpha = 8$ performs best for LLaVA-1V.

For testing on MMMU, we only train on MathVista's disagreements. Whereas for testing on MME, we train on MathVista's and MMMU's disagreements, and analogously for testing on MathVista, we train on MMMU's and MME's disagreements. There are two reasons for this. Firstly, there is an imbalance in the number of disagreements between Mathvista and MME with those of MME outnumbering Mathvista to a large extent. This has repercussions for finetuning as there are a lot more datapoints for MME (yes/no questions) than MathVista (multiple-choice questions). The MMMU test set has multiple-choice questions but a model trained on the MathVista and MME mixture tends to learn the yes/no QA task ignoring multiple-choice questions. As a result, we finetune only on MathVista to learn the multiple-choice task.

B Dataset Statistics

Figure 5 shows pairwise disagreements on MMU and MME datasets for every pair of models in our experiments. Recall that we run debate and consultancy matches only on these disagreements. Table 4 presents the number of reasoning trace samples obtained from debate and consultancy matches.

Dataset	Original	Debate Traces	Consultancy Traces
MME	2,374	4,058	8,120
MMMU	847	2,991	5,978
MathVista	540	1,610	3,220

Table 4: Reasoning traces (total count) obtained through debate and consultancy protocols for our three datasets.

C Licences

823 824 825

810

811

812

814

815

816

818

For all our experiments, we use the versions of the models and datasets released on the HuggingFace platform. All the artifacts but the deepseek-ai's DeepSeek-R1-Distill-Qwen-32B are released under the Apache-2.0 license. The deepseek's model is released under MIT license.

D System Prompts

Judge Prompt

You are a judge in a debate where the participants are arguing over the correct answer to a question grounded in an image. Your task is to determine which answer is correct, based on the arguments provided by the debaters. The arguments contain premises that allow the debaters to come to the answer. The premises are evidence based on the image that only the debaters have access to. However, you haven't seen the image, making your task challenging. To guide you, debaters can only include premises that they observe from the image. They can use their premises to prove their own answer or disprove their opponents answer. Here is the question and the two answers:

Question: <QUESTION>

I:	<answer_a></answer_a>
II:	<answer_b></answer_b>

<NAME_A> is arguing for answer I, and <NAME_B> is arguing for answer II.

Argument Evaluation Guidelines:

1. Relevance Score:

- A premise of an argument should be seen as relevant if it contributes to the acceptance or rejection of the answer the debater is arguing for.
- Rate the premises on a scale of 3 for relevance with the following rubric: i. 3 (High) ii. 2 (Medium) iii. 1 (Low)

2. Acceptability Score:

- A premise of an argument is acceptable if it is worthy of being believed, i.e. you logically think the argument makes sense based on common sense.
- A premise that is inconsistent with any of the premises from the same debater, must receive a low score for acceptability.
- Rate the premises on a scale of 3 for acceptability with the following rubric: i. 3 (High) ii. 2 (Medium) iii. 1 (Low)

3. Credibility Score:

- A premise should be seen as successful in creating credibility if it provides as much detail of the topic of the debate.
- A premise that is based on the description of the image rather than the image that the debaters see, should receive a low score for credibility.
- Rate the premises on a scale of 3 for credibility with the following rubric: i. 3 (High) ii. 2 (Medium) iii. 1 (Low)
- 4. Consistency with Image description:
 - The debaters provided a detailed description apart from the answers to the question. Check if their image description align with their premises. * Note that the premises might contain additional information compared to the image description. However, premises should *not* oppose the description provided by the same debater.
 - Rate the premises on a scale of 3 for consistency with image description with the following rubric: i. 3 (High) ii. 2 (Medium) iii. 1 (Low)

<NAME_A> is arguing for answer I, and <NAME_B> is arguing for answer II. Here's the current transcript of the debate:

<transcript> <TRANSCRIPT> </transcript>

- 1. List down premises from both the debaters.
- 2. For each premise from both the debaters, carry out following steps based on the Argument Evaluation Guidelines:
 - Step 1: Give explanation for its relevance followed by the relevance score.
 - Step 2: Give explanation for its acceptability followed by the acceptability score.
 - Step 3: Give explanation for its credibility followed by the credibility score.
- 3. Filter premises from both the debaters that have received high score of 3 in all the aspects (relevance, acceptability, credibility).

- 4. For each debater check if the Image description provided is consistent with the filtered premises that after providing your explanation for consistency with image description.
 <NAME_A> has provided description of the image as <DESCRIPTION_A>
 <NAME_B> has provided description of the image as <DESCRIPTION_B>
- 5. Using the filtered premises and consistency of debaters with their own image description, provide explanation for which is the likely correct answer to the question "<QUESTION>".

Is it I: "<ANSWER_A>" or II: "<ANSWER_B>".

6. Finally, on a new line mention the answer option letter to the question in the format "Answer: <A|B|C|D|E>" from the choices in the question. * Ensure to only mention the answer letter on the final line.

Debater Prompt (Round 1)

You will act as an expert debater named <NAME> engaging in a debate over answers about an image. You will be shown a question and the answer assigned to you followed by the image. Your task is to find supporting premises from the image in favour of the answer assigned to you while another debater will defend the other answer. A blind judge, who can not see the image, will evaluate you and your opponent solely based on your arguments from the image. Your arguments should be based on premises. Premises are observations based on the image which can allow you to truthfully explain your answer.

Here's the current transcript of the debate:

<transcript> <TRANSCRIPT> </transcript>

This is the start of the Round 1 of the debate. Now it's your turn to look at the image, the question, your answer and form logical arguments using observations (premises) from the image.

- Step 1. List all the premises from the image that are relevant to the question "<QUESTION>".
 - Mention all the visible characteristics (like position in the image, color and size) from the image of the objects discussed in the question.
- Step 2. Using these premises, form logical arguments to support answer "<ANSWER_DEFENDING>".

Follow the above steps to provide your response for the ROUND 1 of the debate as <NAME>.

• Remember to only provide response for Round 1 only as <NAME>. Do NOT provide arguments for your opponent <OPPONENT_NAME>.

Debater Prompt (Round 2)

You will act as an expert debater named <NAME> engaging in a debate over answers about an image. You will be shown a question and the answer assigned to you followed by the image. Your task is to find supporting premises from the image in favour of the answer assigned to you while another debater will defend the other answer. A blind judge, who can not see the image, will evaluate you and your opponent solely based on your arguments from the image. Your arguments should be based on premises. Premises are observations based on the image which can allow you to truthfully explain your answer. (Truthfully explain your answer.

<TRANSCRIPT>

</transcript> This is the start of the Round 2 of the debate. Now it is your turn to look at Round 1 of the transcript, and

provide a rebuttal for Round 2 as <NAME> supporting your answer <ANSWER_DEFENDING>.

- Step 1. From the Round 1 of the transcript mentioned above, identify agreements and disagreements with your opponent <OPPONENT_NAME>.
- Step 2. For each of the disagreements with <OPPONENT_NAME>, provide explanation based on your observations of the image.
 - Use all the visible characteristics (like position in the image, color and size) from the image to argue against your opponent <OPPONENT_NAME>.
- Step 3. Using these premises, form logical arguments to support your answer "<ANSWER_DEFENDING>" to the question "<QUESTION>".

820

 Remember to only provide response for Round 1 only as <NAME>. Do NOT provide arguments for your opponent <OPPONENT_NAME>.

Reasoning Trace Extractor Prompt

Your task is to write an answer that is grounded in observations from an image using logical explanations to a question about the image. However you do not have access the image, making your task challenging. You have access to the judgement of a debate between two debaters who can see the image. The debaters are arguing for the correct answer to the question, by providing premises based on their observations from image. The judge has arrived to the final answer based on the arguments and premises from each debater.

Here is the Question of the debate:

<QUESTION>

The two answers being defended by the debaters A and B are:

- <ANSWER A>
- <ANSWER B>

Here are the two descriptions of the images from Debater A and Debater B:

Description from Debater A: <DESCRIPTION A> Description from Debater B: <DESCRIPTION B>

Here is the judgement by the judge for the debate along with the answers:

<JUDGEMENT>

Based on the above details carry out the following steps in <think></think> tags:

- Step 1. List all the premises that the judge uses to come to the correct answer.
- Step 2. From both the descriptions, list all the observations relevant to the answer provided in the judgement.
- Step 3. Filter all the observations from descriptions and premises from the judgement, that can explain the answer.
- Step 4. Write a detailed answer all the points from Step 3 using the following guidelines:
 - \star The answer should be grounded in the image.
 - ★ The answer should be written as if you can see the image. Use premises, judgements and the descriptions to write the detailed answer along with the logical explanation.
 - ★ Do not mention debaters or the judgement, their descriptions, or the judgement to explain the answer. As you need to answer the question as if you can see the image.

<QUESTION>

Write the answer to the question with the explanation and reasoning as if you can see the image as per the guidelines in </answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></answer></

E Examples of System Output

Example Debate (Molmo-O vs. Molmo-D)

ROUND 1

DEBATER A

As Debater A, I will now provide arguments supporting the answer B (No) based on the premises from the image:

- 1. The image shows a map of the United States with different shades representing hospital beds per 1 million people.
- 2. New Jersey and Georgia are both visible on the map.
- 3. New Jersey is located in the northeastern part of the country, while Georgia is in the southeastern part.

- 4. The color gradient on the map suggests a variation in hospital bed availability across states.
- 5. The darker shades on the map generally indicate higher values, while lighter shades indicate lower values.
- 6. Based on the color distribution, Georgia appears to have a darker shade compared to New Jersey.
- 7. This suggests that Georgia likely has a higher value (more hospital beds per 1 million people) than New Jersey.
- 8. Therefore, the answer B (No) is correct, as Georgia does not have a higher value than New Jersey.

These premises are solely based on the visual information provided in the image and form a logical argument supporting the answer B.

DEBATER B

New Jersey's hospital bed value is higher than Georgia's. The image shows a map of the United States with hospital bed values assigned to each state. New Jersey is colored in a darker shade of red compared to Georgia, which indicates a higher value. Additionally, New Jersey is positioned in the northeastern part of the country, an area known for its higher healthcare needs and population density, which typically correlates with more hospital beds per capita.

ROUND 2

DEBATER A

As Debater A, I will now provide additional arguments supporting the answer B (No) based on the premises from the image:

- 1. The map clearly shows a color gradient representing hospital beds per 1 million people across the United States.
- 2. New Jersey is distinctly colored in a darker shade compared to Georgia.
- 3. The darker color indicates a higher value for hospital beds per 1 million people in New Jersey.
- 4. Georgia, on the other hand, is colored in a lighter shade, suggesting a lower value.
- 5. This color difference is a clear visual indicator of the relative hospital bed availability between the two states.
- 6. The positioning of New Jersey and Georgia on the map is consistent with their actual locations, allowing for easy comparison of their values.
- 7. The use of a color gradient is a common visual representation technique to show differences in data values, making it a reliable indicator in this context.
- 8. The contrast between the two states' colors is quite noticeable, further supporting the conclusion that New Jersey has a higher value than Georgia.

These premises are solely based on the visual information provided in the image and form a logical argument supporting the answer B.

Debater B

New Jersey does not have a higher value than Georgia. The image shows a map of the United States with hospital bed values assigned to each state. New Jersey is colored in a lighter shade of red compared to Georgia, which indicates a lower value. Additionally, New Jersey is positioned in the northeastern part of the country, an area known for its higher healthcare needs and population density, which typically correlates with more hospital beds per capita. However, the color gradient on the map suggests a variation in hospital bed availability across states. The darker shades on the map generally indicate higher values, while lighter shades indicate lower values. Based on the color distribution, Georgia appears to have a darker shade compared to New Jersey. This suggests that Georgia likely has a higher value (more hospital beds per 1 million people) than New Jersey. Therefore, the answer A (Yes) is correct, as New Jersey does not have a higher value than Georgia. These premises are solely based on the visual information provided in the image and form a logical argument supporting the answer A.

Example Judgment			
Alright, let's try to figure out the correct answer to the question: "Is the number of gray cars less than the number of small metallic minivans?" The options are (A) Yes and (B) No.			
First, I'll go throu provided.	gh the arguments from both debaters, A and B, and evaluate their premises based on the guidelines		
Starting with Deba	ater A's premises:		
1. Observation	n: There are two gray cars visible in the image.		
2. Observation	n: There is one small metallic minivan visible in the image.		
3. Conclusion	: The number of gray cars (2) is less than the number of small metallic minivans (1).		
Relevance:	These premises directly address the question by counting the number of gray cars and minivans. So, relevance is high (3).		
Acceptability:	It makes sense to count the number of each vehicle. However, the issue arises when considering whether the tandem bicycle is a car. If it's not, then the count might be incorrect. So, acceptability is medium (2).		
Credibility:	The premises are based on the image, but the description from Debater A doesn't mention the tandem bicycle, which might be a problem. So, credibility is low (1).		
Consistency:	Debater A's image description doesn't mention any vehicles, only clouds and a sky. This inconsistency means the premises aren't consistent with the description (1).		
Now, looking at D	ebater B's premises:		
• Premise 1: 7	There are two gray cars in the image, one is a tandem bicycle and the other is a small metallic minivan.		
• Premise 2: The tandem bicycle is not a car but a vehicle designed for two riders.			
• Premise 3:	The small metallic minivan is a car and there is only one such vehicle in the image.		
Relevance:	These premises are relevant as they clarify what counts as a car. Relevance is high (3).		
Acceptability:	It's logical to distinguish between cars and bicycles. However, the description mentions a motorcycle, which might affect the count. Acceptability is medium (2).		
Credibility:	The premises are based on the image, which includes the tandem bicycle. Credibility is high (3).		
Consistency:	Debater B's image description includes the tandem bicycle, SUVs, and a motorcycle, aligning with the premises. Consistency is high (3).		
After evaluating, only Debater B's premises have high scores in all aspects. The image description from B includes the tandem bicycle, which is gray but not a car. The small metallic minivan is the only car. Therefore, the number of gray cars (1) is equal to the number of minivans (1), so the answer is (B) No.			
The correct provide	nic determined by avaluating the promises and image descriptions. Dehoter D'a promises are consistent		

The correct answer is determined by evaluating the premises and image descriptions. Debater B's premises are consistent with their image description, correctly identifying the tandem bicycle as not a car. Thus, the number of gray cars equals the number of minivans.

Answer: B

Example Trace (for Question and Image in Figure 1)

Yes, the blue umbrella is under the black umbrella. The image shows a boat with passengers holding various umbrellas, including a black and a blue one. The black umbrella is positioned at the front of the boat, while the blue umbrella is directly below it. Due to the angle and positioning, the blue umbrella appears partially obscured by the black umbrella. This arrangement indicates that the blue umbrella is indeed under the black umbrella.

836

F Dataset Examples





MME	
	Is a python code shown in the picture?
<pre>print ("Hello, C++!")</pre>	Please answer yes or no.
	Choices: Yes/No
	Answer: Yes





MMMU



What is NOT exhibited in the painting?			
Choices:	A: hierarchical scaleB: graphic representation of horror and despairC: a sense of immediacy and dramaD: use of sharply contrasting light and shade		
Answer:	А		