

# LOGICAL-COMMONSENSEQA: A Benchmark for Logical Commonsense Reasoning

Anonymous ACL submission

## Abstract

Commonsense reasoning often involves evaluating multiple plausible interpretations rather than selecting a single atomic answer, yet most benchmarks rely on single-label evaluation, obscuring whether statements are jointly plausible, mutually exclusive, or jointly implausible. We introduce LOGICAL-COMMONSENSEQA, a benchmark that reframes commonsense reasoning as *logical composition* over pairs of atomic statements using plausibility-level operators (AND, OR, NEITHER/NOR). Evaluating instruction-tuned, reasoning-specialized, and fine-tuned models under zero-shot, few-shot, and chain-of-thought prompting, we find that while models perform reasonably on conjunctive and moderately on disjunctive reasoning, performance degrades sharply on negation-based questions. LOGICAL-COMMONSENSEQA exposes fundamental reasoning limitations and provides a controlled framework for advancing compositional commonsense reasoning.

## 1 Introduction

Commonsense reasoning is central to human cognition and a long-standing challenge in artificial intelligence and natural language understanding. Developmental studies suggest it emerges from intuitive, experience-based frameworks in early childhood and becomes more abstract and socially grounded through embodied experience, social interaction, cultural transmission, and language (Spelke et al., 1992; Spelke and Kinzler, 2007; Gopnik, 2004; Meltzoff et al., 2009). In contrast, large language models (LLMs) acquire “commonsense-like” knowledge indirectly from large-scale text corpora via objectives such as next-token prediction, rather than through grounded interaction or social learning (Bender and Koller, 2020). Despite this, LLMs exhibit emergent behaviors resembling human intuitive reasoning, including inferring intentions and anticipating socially plausible actions (Sap et al.,

2019a; Bisk et al., 2020). However, this competence is brittle, as models fail under minimal semantic perturbations, struggle with logical consistency, and rely on statistical shortcuts rather than genuine reasoning (Ettinger, 2020; Geiger et al., 2020; McCoy et al., 2019; Niven and Kao, 2019). To improve the evaluation of commonsense reasoning, we introduce LOGICAL-COMMONSENSEQA, a benchmark that reframes the task as a *logical composition over multiple plausible answers*, requiring models to reason about relationships between plausible statements rather than select a single response.

Early approaches to commonsense reasoning relied on symbolic knowledge bases (Lenat et al., 1990; Liu and Singh, 2004), which were expressive but limited in scalability and coverage. More recent benchmarks shift the evaluation towards the ability of neural models to infer plausible answers from text (Talmor et al., 2019; Sap et al., 2019b,a; Zellers et al., 2019; Bisk et al., 2020; Sakaguchi et al., 2021). However, these benchmarks reduce commonsense reasoning to single-answer prediction, despite the fact that many real-world situations admit multiple plausible interpretations (Min et al., 2020). For example, in the question “*The fox walked from the city into the forest—what was it looking for?*” both “*natural habitat*” and “*shelter from disturbances*” are plausible answers. Restricting the evaluation to a single atomic choice removes this natural ambiguity and obscures whether models recognize jointly plausible, mutually exclusive, or jointly implausible interpretations.

Beyond linguistic ambiguity, commonsense is both socially grounded and compositional. People judge commonsense statements not only by personal belief but also by estimating whether others would share that belief (Whiting and Watts, 2024), implying that commonsense reflects a distribution over plausible interpretations rather than a single canonical truth. At the same time, human reasoning integrates multiple alternatives through composi-

tional inference, recognizing when statements are jointly plausible (AND), partially plausible (OR), or jointly implausible (NEITHER/NOR). These relational judgments are central to everyday reasoning about intentions, plans, and social norms, yet current benchmarks fail to evaluate whether models can reason over such socially grounded and compositional plausibility relationships.

In LOGICAL-COMMONSENSEQA, each instance pairs two atomic answers under one of three symbolic relations—AND, OR, or NEITHER/NOR. These relations are not strict propositional operators, but *plausibility-level composition constructs* indicating joint, partial, or joint implausibility. This formulation preserves the multiple-choice question (MCQ) format of modern benchmarks while explicitly modeling ambiguity and compositional reasoning. We evaluate instruction-tuned, reasoning-specialized, and fine-tuned LLMs under zero-shot, few-shot, and chain-of-thought prompting. Results show that while models perform reasonably well on conjunctive and moderately on disjunctive reasoning, performance degrades sharply on negation-based compositions, revealing persistent gaps in the commonsense reasoning capabilities of state-of-the-art systems.

## 2 Related Work

Commonsense reasoning has been widely studied through symbolic knowledge bases, neural models, and benchmark-driven evaluation. Early symbolic approaches such as CYC (Lenat et al., 1990) and CONCEPTNET (Liu and Singh, 2004) enabled structured inference but suffered from limited coverage and scalability. More recent benchmarks, including COMMONSENSEQA (Talmor et al., 2019), SOCIALIQA (Sap et al., 2019b), ATOMIC (Sap et al., 2019a), PIQA (Bisk et al., 2020), HELLASWAG (Zellers et al., 2019), and WINOGRANDE (Sakaguchi et al., 2021), shift evaluation towards the ability of neural models to infer plausible responses from text. However, these benchmarks adopt a single atomic-answer format that fails to capture the inherent ambiguity of many real-world commonsense questions.

Prior work has emphasized that ambiguity and human disagreement are intrinsic to commonsense reasoning. AMBIGQA (Min et al., 2020) and PROTOQA (Boratko et al., 2020) show that many questions naturally admit multiple valid answers, motivating the need to evaluate systems beyond sin-

gle gold answers. Separately, logical reasoning benchmarks such as LOGIQA (Liu et al., 2020), RECLOR (Yu et al., 2020), and COM2 (Fang et al., 2024) evaluate deductive inference using formal logical operators, but target logical validity rather than commonsense plausibility. Complementary work in the social sciences argues that commonsense is socially grounded, reflecting shared expectations rather than objective truth (Whiting and Watts, 2024).

Our work differs from these lines of research by introducing LOGICAL-COMMONSENSEQA, a benchmark that evaluates commonsense reasoning as *compositional plausibility*. Instead of open-ended multi-answer annotation or formal deductive logic, LOGICAL-COMMONSENSEQA preserves the MCQ format of modern commonsense benchmarks while explicitly encoding joint plausibility, partial plausibility, and joint implausibility through symbolic composition over human-validated atomic answers.

## 3 The Logical Commonsense Dataset

We construct LOGICAL-COMMONSENSEQA to evaluate compositional, multi-answer commonsense reasoning by extending COMMONSENSEQA (Talmor et al., 2019), a benchmark of 12,247 MCQs with one correct answer and four distractors that probes commonsense relations such as cause-effect, purpose, and spatial association. LOGICAL-COMMONSENSEQA reformulates these single-answer questions into logically composed items that explicitly model relationships among plausible alternatives while preserving full compatibility with the MCQ format.

**Construction Pipeline.** The dataset is built using a three-stage pipeline that integrates neural generation with deterministic symbolic composition, converting atomic answers into logically structured multi-answer instances.

*Stage 1: Generation of Candidate Options.* Starting from COMMONSENSEQA questions and their gold answers, we sample 5,000 instances and use GPT-4o-mini to over-generate a diverse set of atomic answer candidates. The model is prompted to produce both plausible and implausible alternatives, with an emphasis on multi-step causal or situational reasoning rather than shallow lexical cues. This process yields 4-6 plausible and implausible candidate answers per question, spanning physical, social, and situational commonsense domains.

Operator	Interpretation	Example
$a$ AND $b$	Both statements independently plausible	<i>Q: Sammy wanted to go to where the people were. Where might he go?</i> A: local events AND social venues
$a$ OR $b$	At least one statement plausible	<i>Q: Sammy wanted to go to where the people were. Where might he go?</i> A: local events OR empty parks
NEITHER $a$ NOR $b$	Neither statement plausible	<i>Q: Sammy wanted to go to where the people were. Where might he go?</i> A: NEITHER quiet retreats NOR empty parks

Table 1: Composition operators used in LOGICAL-COMMONSENSEQA to express plausibility relations.

*Stage 2: Refinement and Pruning.* We then use GPT-4o-mini to refine and filter the generated candidates according to three criteria: (1) removing logically inconsistent or factually incorrect answers; (2) eliminating trivial options that can be resolved through keyword matching; and (3) identifying highly plausible options that satisfy most semantic and contextual constraints of the question but fail due to a non-obvious commonsense violation. This stage yields a curated set of three correct and four incorrect atomic options per question, each requiring multi-step inference.

*Stage 3: Deterministic Logical Composition.* In the final stage, we use a symbolic program to deterministically combine pairs of refined atomic options into logical compositions labeled with one of the three relations shown in Table 1. These relations are not strict propositional operators; instead, they encode commonsense judgments about joint plausibility. This procedure yields 14,997 instances in total, with 4,999 instances per base relation type.

The prompts for all stages of the construction pipeline are provided in Appendix A.3.

**Human Validation** To ensure the commonsense validity of the refined atomic options, we conducted a human validation study using the awareness–consensus framework proposed by Whiting and Watts (2024). Two independent annotators evaluated a random sample of 250 questions (5% of Stage 2). For each atomic option, annotators assessed (1) whether they personally believed the option was correct (awareness) and (2) whether they believed most others would agree (consensus). This dual assessment captures both individual plausibility and perceived social agreement, aligning with recent views of commonsense as socially grounded. Disagreements between annotators were resolved through adjudication.

We report label accuracy on the adjudicated subset covering 50% of the question-level test split defined after Stage 2. Performance is high for AND (89.2%), OR (96.4%), and MIXED (88.4%), and

reasonable for NEITHER/NOR (73.6%). These numbers reflect a reasonably accurate labeling of composed instances.

Inter-annotator agreement, measured using Cohen’s  $\kappa$ , indicated moderate agreement on human judgments ( $\kappa = 0.49$ ) before adjudication, suggesting that the refined options are both semantically coherent and socially interpretable.

**Resulting Dataset** After Stage 2, we obtain 4,999 curated questions with refined atomic answer sets. In Stage 3, each question is expanded into logically structured instances evenly distributed across the three base relation types, with an additional MIXED setting in which operators are randomly assigned across options. After filtering, the final dataset contains 19,996 instances, with AND, OR, NEITHER/NOR, and MIXED each contributing 4,999 instances.

## 4 Task Formulation

Instances in LOGICAL-COMMONSENSEQA consist of a question in natural language and four options. Unlike standard benchmarks where each option is an atomic answer, each option in our dataset is a *logical composition* of two independent atomic statements combined using one of three plausibility-level operators: AND, OR, or NEITHER/NOR (see Table 1).

Given a question, the model must select the option whose composite plausibility best matches the implied commonsense constraints. While the task remains a simple MCQ problem, it requires models to assess the plausibility of individual statements and reason about their interaction under the specified operator. We include both operator-specific conditions, in which all answer options for a question share the same logical operator, and a MIXED condition, in which different operators appear across options. While the operator-specific conditions allow for controlled evaluation of each relation type, the MIXED condition prevents models from exploiting operator-specific patterns and

requires direct inference over the composed statements. Overall, this formulation preserves the MCQ format while substantially increasing the demands on compositional commonsense reasoning.

## 5 Experiments

We evaluate a diverse set of encoder-only classifiers, encoder–decoder models, and decoder-only LLMs on LOGICAL-COMMONSENSEQA under zero-shot, few-shot, and supervised fine-tuning regimes. Full details on model selection, prompting strategies, and training and inference procedures are provided in Appendices A.1, A.2, and A.5.

Following the dataset construction described in Section 3, we use all 19,996 instances of LOGICAL-COMMONSENSEQA split into 11,996 training, 6,000 development, and 2,000 test examples. Splits are stratified by logical relation, ensuring an even 25% distribution across AND, OR, NEITHER/NOR, and MIXED conditions. We report prompt accuracy and macro  $F_1$  as evaluation metrics.

For analysis, we further divide the test set into a human-validated (HV) subset and a non-validated (NV) subset, each comprising 50% of the test data. Table 2 summarizes the main results. Across all model families, two consistent patterns emerge: (1) all models perform reasonably well on AND and moderately on OR, and (2) performance collapses on NEITHER/NOR in zero- and few-shot settings, even for the strongest LLMs. Observed weaknesses are amplified in the MIXED condition, where different operators appear across options and models must implicitly infer the governing relation.  $F_1$  drops to 43–56% for all few-shot models, suggesting that these models do not reliably maintain operator-level representations and instead revert to surface-level heuristics. In contrast, fine-tuned models achieve 83–93%  $F_1$  across all operators, suggesting that the task is learnable with supervision and that zero- and few-shot failures, particularly on negation, reflect inference-time limitations rather than dataset artifacts.

Additional results—including varying numbers of shots, prompting strategies, and performance on the NV split—are reported in Appendix A.6. A comprehensive error analysis is also provided in Appendix A.7.

*Comparison with CommonsenseQA.* Finally, Table 3 underscores the central motivation for our benchmark. While LLaMA-3.1-8B achieves 72.2% accuracy on COMMONSENSEQA, its performance

Paradigm	Model	AND	OR	NN	MIX
0-Shot	LLaMA-3.3-70B	80.9	70.9	13.4	53.0
	LLaMA-3.1-8B	71.9	62.2	13.1	41.8
	Qwen2.5-7B	79.6	68.9	12.9	53.2
3-Shot	LLaMA-3.3-70B	85.7	77.5	10.7	50.4
	LLaMA-3.1-8B	71.3	61.3	6.1	42.8
	Qwen2.5-7B	77.9	71.9	6.8	51.4
Fine-tuned	Flan-T5-base	92.8	92.4	89.2	89.6
	DeBERTa-v3-base	87.6	87.2	84.8	82.4

Table 2: **Macro- $F_1$  on the human-validated test set** across plausibility relations. NN = NEITHER/NOR, MIX = MIXED

Dataset	Acc	$F_1$
<b>CommonsenseQA</b>		
LLaMA-3.1-8B	72.2	72.2
<b>LOGICAL-COMMONSENSEQA</b>		
AND	72.0	71.9
OR	62.2	62.2
NEITHER/NOR	13.9	13.1
MIXED	42.7	41.8

Table 3: 0-Shot performance of LLaMA-3.1-8B on the original COMMONSENSEQA versus each logical operator in LOGICAL-COMMONSENSEQA.

drops sharply on LOGICAL-COMMONSENSEQA: 72% on AND, 62.2% on OR, 42.7% on MIXED, and only 13.9% on NEITHER/NOR. This discrepancy shows that benchmarks like COMMONSENSEQA substantially overestimate models’ commonsense reasoning ability by failing to test relational and compositional plausibility judgments.

## 6 Conclusions and Future Work

We introduce LOGICAL-COMMONSENSEQA, a benchmark that reframes commonsense reasoning as selecting among logically composed answer options. By combining human-validated atomic statements with symbolic operators that encode conjunction, disjunction, and negation of plausibility, our work exposes reasoning failures obscured by traditional single-answer benchmarks. Experiments reveal that while LLMs handle conjunctive and, to a lesser extent, disjunctive reasoning, they fail sharply on negation-based composition, highlighting persistent gaps in plausibility evaluation and operator grounding. Future work includes extending our resource to generative settings and structured reasoning frameworks, expanding the operator set, scaling awareness-consensus validation, and studying transfer to open-ended QA, dialog, and planning tasks.

## 343 Limitations

344 LOGICAL-COMMONSENSEQA offers a controlled  
345 framework for probing compositional common-  
346 sense reasoning, but several limitations remain.  
347 The logical operators (AND, OR and NEI-  
348 THER/NOR) capture plausibility-level relations  
349 rather than full propositional semantics, and do not  
350 cover richer structures such as implication, exclu-  
351 sivity or temporal and causal reasoning. Extending  
352 the operator space would enable a more complete  
353 evaluation of symbolic and relational inference.

354 Model coverage is another limitation. Decoder-  
355 only LLMs are evaluated only via prompting, while  
356 smaller encoder-based models receive supervised  
357 training. A broader comparison across instruction-  
358 tuned and specialized architectures may clarify  
359 which components most improve logical composi-  
360 tion.

## 361 Ethical Considerations

362 LOGICAL-COMMONSENSEQA is derived from the  
363 publicly available COMMONSENSEQA dataset and  
364 consists of generic, hypothetical scenarios that do  
365 not reference identifiable individuals or personal  
366 contexts. To the best of our knowledge, the dataset  
367 does not contain private or sensitive information.

368 We used an LLM to initially generate atomic  
369 answer options and increase coverage and diver-  
370 sity. These candidates are subsequently refined and  
371 partially validated by human annotators using an  
372 awareness–consensus framework. As a result, the  
373 dataset reflects socially grounded plausibility judg-  
374 ments rather than objective ground truth, and may  
375 thus inherit cultural or contextual biases.

## 376 References

377 Emily M. Bender and Alexander Koller. 2020. [Climbing](#)  
378 [towards NLU: On meaning, form, and understanding](#)  
379 [in the age of data](#). In *Proceedings of the 58th Annual*  
380 *Meeting of the Association for Computational Lin-*  
381 *guistics*, pages 5185–5198, Online. Association for  
382 Computational Linguistics.

383 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi,  
384 and 1 others. 2020. Piqa: Reasoning about physical  
385 commonsense in natural language. In *Proceedings*  
386 *of the AAAI conference on artificial intelligence*, vol-  
387 *ume 34*, pages 7432–7439.

388 Michael Boratko, Xiang Li, Tim O’Gorman, Rajarshi  
389 Das, Dan Le, and Andrew McCallum. 2020. Pro-  
390 toqa: A question answering dataset for prototypi-  
391 cal common-sense reasoning. In *Proceedings of the*

*2020 Conference on Empirical Methods in Natural*  
*Language Processing (EMNLP)*, pages 1122–1136.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret  
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi  
Wang, Mostafa Dehghani, Siddhartha Brahma, and  
1 others. 2024. [Scaling instruction-finetuned lan-](#)  
[guage models](#). *Journal of Machine Learning Re-*  
*search*, 25(70):1–53.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,  
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,  
Akhil Mathur, Alan Schelten, Amy Yang, Angela  
Fan, and 1 others. 2024. The llama 3 herd of models.  
*arXiv e-prints*, pages arXiv–2407.

Allyson Ettinger. 2020. What bert is not: Lessons from  
a new suite of psycholinguistic diagnostics for lan-  
guage models. *Transactions of the Association for*  
*Computational Linguistics*, 8:34–48.

Tianqing Fang, Zeming Chen, Yangqiu Song, and An-  
toine Bosselut. 2024. Complex reasoning over logi-  
cal queries on commonsense knowledge graphs.  
*arXiv preprint arXiv:2403.07398*.

Atticus Geiger, Kyle Richardson, and Christopher Potts.  
2020. Neural natural language inference models par-  
tially embed theories of lexical entailment and nega-  
tion. In *Proceedings of the Third BlackboxNLP Work-*  
*shop on Analyzing and Interpreting Neural Networks*  
*for NLP*, pages 163–173.

Alison Gopnik. 2004. Finding our inner scientist.  
*Daedalus*, 133(1):21–28.

Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021.  
[Debertav3: Improving deberta using electra-style pre-](#)  
[training with gradient-disentangled embedding shar-](#)  
[ing](#). *arXiv preprint arXiv:2111.09543*.

Douglas B Lenat, Ramanathan V. Guha, Karen Pittman,  
Dexter Pratt, and Mary Shepherd. 1990. Cyc: toward  
programs with common sense. *Communications of*  
*the ACM*, 33(8):30–49.

Hugo Liu and Push Singh. 2004. Conceptnet—a practi-  
cal commonsense reasoning tool-kit. *BT technology*  
*journal*, 22(4):211–226.

Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang,  
Yile Wang, and Yue Zhang. 2020. Logiqa: A  
challenge dataset for machine reading compre-  
hension with logical reasoning. *arXiv preprint*  
*arXiv:2007.08124*.

R. Thomas McCoy, Ellie Pavlick, and Tal Linzen. 2019.  
[Right for the wrong reasons: Diagnosing syntactic](#)  
[heuristics in natural language inference](#). In *Proceed-*  
*ings of the 57th Annual Meeting of the Association for*  
*Computational Linguistics*, pages 3428–3448, Flo-  
rence, Italy. Association for Computational Linguis-  
tics.

444	Andrew N Meltzoff, Patricia K Kuhl, Javier Movellan, and Terrence J Sejnowski. 2009. Foundations for a new science of learning. <i>science</i> , 325(5938):284–288.	<i>Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4149–4158.	499
445			500
446			501
447			
448	Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. <b>AmbigQA: Answering ambiguous open-domain questions</b> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5783–5797, Online. Association for Computational Linguistics.	Mark E Whiting and Duncan J Watts. 2024. A framework for quantifying individual and collective common sense. <i>Proceedings of the National Academy of Sciences</i> , 121(4):e2309535121.	502
449			503
450			504
451			505
452		Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. Reclor: A reading comprehension dataset requiring logical reasoning. <i>arXiv preprint arXiv:2002.04326</i> .	506
453			507
454			508
455	Timothy Niven and Hung-Yu Kao. 2019. <b>Probing neural network comprehension of natural language arguments</b> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4658–4664, Florence, Italy. Association for Computational Linguistics.	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. <b>HellaSwag: Can a machine really finish your sentence?</b> In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4791–4800, Florence, Italy. Association for Computational Linguistics.	510
456			511
457			512
458			513
459			514
460			515
461	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, and 25 others. 2025. <b>Qwen2.5 technical report</b> . <i>Preprint</i> , arXiv:2412.15115.		
462			
463			
464			
465			
466			
467			
468	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106.		
469			
470			
471			
472	Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019a. Atomic: An atlas of machine commonsense for if-then reasoning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 33, pages 3027–3035.		
473			
474			
475			
476			
477			
478			
479	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. 2019b. Socialliqa: Commonsense reasoning about social interactions. In <i>Conference on Empirical Methods in Natural Language Processing</i> .		
480			
481			
482			
483			
484	Elizabeth S Spelke, Karen Breinlinger, Janet Macomber, and Kristen Jacobson. 1992. Origins of knowledge. <i>Psychological review</i> , 99(4):605.		
485			
486			
487	Elizabeth S Spelke and Katherine D Kinzler. 2007. Core knowledge. <i>Developmental science</i> , 10(1):89–96.		
488			
489	Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2022. Entailer: Answering questions with faithful and truthful chains of reasoning. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2078–2093.		
490			
491			
492			
493			
494	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for</i>		
495			
496			
497			
498			
		<b>A Appendix</b>	516
		<b>A.1 Model Selection</b>	517
		This section describes the architectures and evaluation settings of all models considered in our experiments, covering encoder-only, encoder-decoder, and decoder-only LLMs.	518
			519
			520
			521
		<b>Encoder Models</b> We fine-tune DEBERTA-V3-BASE (He et al., 2021), a pretrained encoder which encodes the concatenation of the question and one composite option and applies a classification head to score the four candidate answers.	522
			523
			524
			525
			526
		<b>Encoder-Decoder and Multi-Step Reasoning Models</b> We also fine-tune sequence-to-sequence models to generate the output label. FLAN-T5-BASE (Chung et al., 2024) is fine-tuned to output a single token corresponding to the chosen option. Additionally, we evaluate ENTAILER-11B (Tafjord et al., 2022), a model designed to produce faithful chains-of-reasoning by recursively generating premises through backward-chaining and verifying them via self-querying. Although originally developed for faithful explanation generation, we use Entailer as a strong encoder-decoder baseline to examine whether reasoning-oriented models better capture the compositional structures in our dataset.	527
			528
			529
			530
			531
			532
			533
			534
			535
			536
			537
			538
			539
			540
		<b>Decoder-Only LLMs</b> We evaluate several open-weight instruction-tuned LLMs in zero- and few-shot settings: LLAMA-3.3-70B-INSTRUCT, LLAMA-3.1-8B-INSTRUCT (Dubey et al., 2024) and QWEN2.5-7B-INSTRUCT (Qwen et al., 2025). These models are not fine-tuned; instead, they are prompted to output a single capital letter corresponding to the predicted option. All models are	541
			542
			543
			544
			545
			546
			547
			548

549 evaluated on each logical subset defined in Table 1  
550 as well as on the Mixed setting.

## 551 A.2 Prompting Setup

552 In this section, we describe the prompting strategy  
553 used for decoder-only models. For decoder-only  
554 LLMs, we adopt a unified multiple-choice prompt-  
555 ing scheme. Each prompt presents the question,  
556 the four composite options and an instruction to  
557 answer with a single letter (A-D). We evaluate 0-,  
558 1-, 2- and 3-shot settings using a fixed set of labeled  
559 examples.

560 Decoding uses temperature = 0.0, top-p =  
561 0.9, and max\_new\_tokens= 3. Outputs not corre-  
562 sponding to a valid answer option are treated as  
563 incorrect.

564 **Chain-of-Thought (CoT)** We evaluate CoT  
565 prompting for LLAMA-3.1-8B-INSTRUCT only,  
566 due to computational cost. CoT runs use a larger  
567 generation budget (max\_new\_tokens = 200-300).

## 568 A.3 Prompt Templates

569 This section provides the exact prompt templates  
570 used in our experiments. Prompts differ only in  
571 the number of in-context examples and the logical  
572 operator expressed in the answer options.

573 **Zero-shot Prompt** In the zero-shot setting, we  
574 don't provide any labeled examples.

575 Answer the following commonsense question by  
576 selecting the correct option. Respond with **only**  
577 the single capital letter (A, B, C, or D).

578 Question: <QUESTION>

- 579 A. <OPTION A>
- 580 B. <OPTION B>
- 581 C. <OPTION C>
- 582 D. <OPTION D>

583  
584 Answer:

585 **Few-shot Prompting** For few-shot (1-, 2-, and  
586 3-shot) evaluation, we prepend labeled examples  
587 before the target question. **The 1- and 2-shot**  
588 **prompts are strict prefixes of the 3-shot prompt,**  
589 using the first one or two examples respectively.

590 For operator-specific subsets (AND, OR, NEI-  
591 THER/NOR), all in-context examples match the  
592 operator being evaluated. For the MIXED setting,  
593 examples are sampled across operators.

## 594 Conjunctive Reasoning (AND) Prompt

595 Question: Sammy wanted to go to where the peo-  
596 ple were. Where might he go?

- A. social venues AND quiet retreats 597
- B. local events AND social venues 598
- C. sports arenas AND quiet retreats 599
- D. sports arenas AND train platforms 600

601  
602 Answer: B

603 Question: The fox walked from the city into the  
604 forest, what was it looking for?

- A. suitable prey AND shelter from disturbances 605
- B. urban garden AND farm fields 606
- C. natural habitat AND farm fields 607
- D. suitable prey AND neighborhood pets 608

609  
610 Answer: A

611 Question: What home entertainment equipment  
612 requires cable?

- A. wireless speaker AND portable projector 613
- B. television AND home theater system 614
- C. wireless speaker AND gaming console 615
- D. digital frame AND portable projector 616

617  
618 Answer: B

619 Now answer the following commonsense question  
620 by selecting the correct option. Respond with  
621 **only** the single capital letter (A, B, C, or D).

622 Question: <TARGET AND QUESTION>

- A. <OPTION A> 623
- B. <OPTION B> 624
- C. <OPTION C> 625
- D. <OPTION D> 626

627  
628 Answer:

## 629 Disjunctive Reasoning (OR) Prompt

630 Question: Sammy wanted to go to where the peo-  
631 ple were. Where might he go?

- A. quiet retreats OR empty parks 632
- B. local events OR empty parks 633
- C. sports arenas OR empty parks 634
- D. train platforms OR empty parks 635

636  
637 Answer: B

638 Question: The forgotten leftovers had gotten quite  
639 old; he found it covered in mold in the back of his  
640 what?

- A. living room OR dining area 641
- B. utility drawer OR kitchen counter 642
- C. living room OR utility drawer 643
- D. cold storage OR utility drawer 644

645  
646 Answer: D

647 Question: What do people use to absorb extra ink  
648 from a fountain pen?

- A. absorbent paper OR notebook cover 649
- B. ink storage OR writing surface 650
- C. notebook cover OR writing surface 651
- D. ink storage OR fabric swatch 652

653

654	Answer: A	A. natural habitat AND shelter from disturbances	709
		B. farm fields OR neighborhood pets	710
655	Now answer the following commonsense question	C. city park AND neighborhood pets	711
656	by selecting the correct option. Respond with	D. farm fields AND neighborhood pets	712
657	<b>only</b> the single capital letter (A, B, C, or D).		713
658	Question: <TARGET OR QUESTION>	Answer: A	714
659	A. <OPTION A>	Question: Google Maps and other highway and	715
660	B. <OPTION B>	street GPS services have replaced what?	716
661	C. <OPTION C>	A. historical navigation logs AND tourist	717
662	D. <OPTION D>	information centers	718
663		B. NEITHER printed road maps NOR route	719
664	Answer:	planning software	720
		C. manual navigation techniques AND tourist	721
665	<b>Negation-Based Reasoning (NEITHER/NOR)</b>	information centers	722
666	<b>Prompt</b>	D. manual navigation techniques AND route	723
		planning software	724
			725
667	Question: Sammy wanted to go to where the peo-	Answer: D	726
668	ple were. Where might he go?		
669	A. NEITHER local events NOR train platforms	Question: Where is a business restaurant likely to	727
670	B. NEITHER sports arenas NOR train platforms	be located?	728
671	C. NEITHER social venues NOR sports arenas	A. corporate office clusters AND at university	729
672	D. NEITHER social venues NOR quiet retreats	campuses	730
673		B. business park entrances AND near conference	731
674	Answer: B	venues	732
		C. NEITHER near conference venues NOR at	733
675	Question: The only baggage the woman checked	university campuses	734
676	was a drawstring bag. Where was she heading	D. NEITHER near conference venues NOR near	735
677	with it?	tourist sites	736
678	A. NEITHER family gathering NOR local		737
679	gathering	Answer: B	738
680	B. NEITHER airport travel NOR local gathering		
681	C. NEITHER local gathering NOR short trip	Now answer the following commonsense question	739
682	D. NEITHER business engagement NOR local	by selecting the correct option. Respond with	740
683	gathering	<b>only</b> the single capital letter (A, B, C, or D).	741
684		Question: <TARGET MIXED QUESTION>	742
685	Answer: C	A. <OPTION A>	743
		B. <OPTION B>	744
686	Question: The forgotten leftovers had gotten quite	C. <OPTION C>	745
687	old, he found it covered in mold in the back of his	D. <OPTION D>	746
688	what?		747
689	A. NEITHER cold storage NOR dining area	Answer:	748
690	B. NEITHER sealed container NOR living room		
691	C. NEITHER food cabinet NOR utility drawer	<b>Chain-of-Thought Prompting</b> For CoT evalua-	749
692	D. NEITHER living room NOR utility drawer	tion, we add a brief instruction asking the model to	750
693		reason about the plausibility of each atomic state-	751
694	Answer: D	ment and how the logical operator applies before	752
		providing the final answer.	753
695	Now answer the following commonsense question	Carefully reason about each atomic statement and	754
696	by selecting the correct option. Respond with	the logical operator. Then provide the final answer	755
697	<b>only</b> the single capital letter (A, B, C, or D).	as a single capital letter (A–D).	756
698	Question: <TARGET NEITHER/NOR QUES-		
699	TION>	<b>Decoding Constraints</b> All decoder-only mod-	757
700	A. <OPTION A>	els are evaluated with near-deterministic decoding	758
701	B. <OPTION B>	(temperature=0.0, top-p=0.9). Outputs are con-	759
702	C. <OPTION C>	strained to a single capital letter (A–D); any other	760
703	D. <OPTION D>	output is treated as incorrect.	761
704			
705	Answer:	<b>A.4 Annotation Guidelines</b>	762
706	<b>Mixed Prompt</b>	This section describes the guidelines provided to	763
		human annotators for validating refined atomic an-	764
707	Question: The fox walked from the city into the	swer options in Stage 2 of dataset construction.	765
708	forest, what was it looking for?		

766	<b>Annotation Task</b>	Each annotation item consists of a commonsense question and a set of automatically labeled answer options. Options are pre-labeled as <i>correct</i> or <i>incorrect</i> by the refinement pipeline. Annotators are instructed to evaluate each option <i>independently</i> , without reconsidering or changing the question itself.	812
767			813
768			814
769			815
770			816
771			817
772			
773		The goal of annotation is only <i>verification</i> : annotators judge whether the assigned labels are reasonable under commonsense reasoning.	818
774			819
775			820
776	<b>Evaluation Criteria</b>	For each option, annotators provide judgments along two dimensions, following the awareness–consensus framework of <a href="#">Whiting and Watts (2024)</a> :	821
777			822
778			823
779			824
780	1. <b>Correctness (Awareness).</b>	Annotators indicate whether they personally believe the option is a valid answer to the question.	825
781			826
782			827
783		• For options labeled as <i>correct</i> : “Do you believe this option is actually a valid answer?”	828
784			829
785			830
786		• For options labeled as <i>incorrect</i> : “Do you believe this option is actually wrong or implausible?”	831
787			832
788			833
789	2. <b>Commonsense (Consensus).</b>	Annotators indicate whether they believe most adults would agree with the judgment using everyday knowledge.	834
790			835
791			836
792			837
793		Each criterion is rated using a three-point scale: Yes, No, or Maybe. Annotators are required to provide a brief comment for any No or Maybe response (e.g., “could be true in some contexts” or “requires specialized knowledge”).	838
794			839
795			840
796			841
797			842
798	<b>Definition of Commonsense</b>	Annotators are instructed to interpret commonsense as knowledge that an average adult could reasonably be expected to possess, without relying on formal education, technical expertise, or domain-specific facts. Options requiring specialized scientific, historical, or technical knowledge are marked accordingly.	843
799			844
800			845
801			846
802			847
803			848
804			849
805	<b>Handling Ambiguity</b>	Annotators are encouraged to consider realistic context and pragmatic reasoning. When an option is plausible only under rare or contrived interpretations, it should be marked as Maybe rather than Yes. The goal is to assess typical human judgments, not theoretical possibility.	850
806			851
807			852
808			853
809			854
810			855
811			856
			857
			858
			859
	<b>Adjudication</b>	Two annotators independently evaluate each item. Disagreements, as well as any option receiving a Maybe are resolved through adjudication by a third reviewer. Only disputed cases are adjudicated; non-disputed options retain their original labels.	812
			813
			814
			815
			816
			817
	<b>Illustrative Examples</b>	Annotators are provided with worked examples illustrating borderline cases, such as partially plausible distractors, context-dependent interpretations, and options that are topically related but logically invalid. These examples are used to clarify expectations but are not treated as templates.	818
			819
			820
			821
			822
			823
			824
	<b>A.5 Training and Inference Details</b>		825
		This section reports the training configurations and inference settings used across all evaluated models.	826
		Encoder-only and encoder-decoder models are fine-tuned with AdamW ( $5 \times 10^{-5}$ learning rate, weight decay 0.01) for up to 10 epochs with early stopping (patience 3), using a maximum sequence length of 512 and an effective batch size of 8.	827
		Decoder-only LLMs are evaluated without gradient updates under 0–3 shot prompting (with/without CoT) and outputs are constrained to a single label (A–D). Most experiments run on NVIDIA P100 GPUs. Only the largest models and CoT evaluations are run on A100 GPUs.	828
			829
			830
			831
			832
			833
			834
			835
			836
			837
			838
	<b>A.6 Additional Results</b>		839
		This section reports supplementary experimental results that complement the main findings discussed in Section 5. Table 4 reports accuracy and macro-F1 across all logical relations for human-validated (HV) and non-human-validated (NV) test subsets, including instruction-tuned and fine-tuned models. Table 5 presents macro-F1 for decoder-only LLMs under 0–3 shot prompting and CoT settings. Encoder-only and encoder–decoder models without fine-tuning exhibit uniformly low performance. Overall trends mirror the main paper: reasonable performance on AND and OR, and collapse on NEITHER/NOR across settings.	840
			841
			842
			843
			844
			845
			846
			847
			848
			849
			850
			851
			852
	<b>A.7 Error Analysis</b>		853
		We analyze errors made by LLAMA-3.1-8B in the zero-shot setting. Table 6 summarizes representative failure patterns revealed by LOGICAL-COMMONSENSEQA. Because each option is a compositional pair of atomic statements, errors expose not only mis-classifications but also	854
			855
			856
			857
			858
			859

Model	AND				OR				NN				MIX			
	NV A	NV F1	HV A	HV F1	NV A	NV F1	HV A	HV F1	NV A	NV F1	HV A	HV F1	NV A	NV F1	HV A	HV F1
LLaMA-70B	76.4	76.2	80.8	80.9	72.4	72.3	71.2	70.9	13.6	13.4	14.0	13.4	53.7	53.4	53.2	53.0
LLaMA-8B	67.2	66.7	72.0	71.9	63.2	62.6	62.2	62.2	8.6	7.9	13.9	13.1	44.5	43.9	42.7	41.8
Qwen2.5-7B	76.2	76.2	79.6	79.6	70.4	70.5	68.8	68.9	13.6	11.8	13.3	12.9	54.0	53.7	53.2	53.2
Entailer-11B	5.6	5.5	3.2	3.2	20.8	20.8	18.4	18.4	22.0	22.0	24.0	24.1	16.8	16.8	22.4	22.3
Flan-T5-base	38.8	38.0	37.2	36.2	63.2	62.2	65.2	64.3	12.0	11.4	16.8	16.3	32.8	32.5	33.6	33.0
Flan-T5-base (FT)	94.8	94.8	92.8	92.8	97.6	97.6	92.4	92.4	93.2	93.2	89.2	89.2	87.2	87.2	89.6	89.6
DeBERTa-v3-base	28.8	28.6	27.2	27.2	23.6	23.8	23.2	23.1	30.8	30.9	27.2	27.4	28.8	28.8	28.4	28.4
DeBERTa-v3-base (FT)	90.8	90.8	87.6	87.6	88.0	88.0	87.2	87.2	89.2	89.2	84.8	84.8	84.8	84.8	82.4	82.4

Table 4: **0-shot vs. fine-tuned:** Accuracy (A) and macro-F1 across four logical relations. NV = non-human-validated labels; HV = human-validated subset. NN = NEITHER/NOR, MIX = MIXED, FT = fine-tuned.

Model	AND		OR		NN		MIX	
	NV F1	HV F1	NV F1	HV F1	NV F1	HV F1	NV F1	HV F1
<b>0-shot</b>								
LLaMA-70B	76.2	80.9	72.3	70.9	13.4	13.4	53.4	53.0
LLaMA-8B	66.7	71.9	62.6	62.2	7.9	13.1	43.9	41.8
Qwen-7B	76.2	79.6	70.5	68.9	11.8	12.9	53.7	53.2
<b>1-shot</b>								
LLaMA-70B	80.3	84.8	76.9	79.5	7.9	11.3	54.7	53.1
LLaMA-8B	72.1	67.0	61.7	62.0	5.9	5.3	45.0	41.7
Qwen-7B	73.9	79.6	75.0	71.6	5.7	7.4	57.4	55.4
<b>2-shot</b>								
LLaMA-70B	78.7	83.6	74.0	76.8	9.1	12.2	53.1	51.7
LLaMA-8B	71.4	75.4	61.6	63.9	6.0	6.6	46.1	41.2
Qwen-7B	72.3	77.5	70.6	69.6	5.7	6.9	55.2	52.0
<b>3-shot</b>								
LLaMA-70B	80.0	85.7	75.1	77.5	7.8	10.7	52.4	50.4
LLaMA-8B	75.8	71.3	61.0	61.3	5.5	6.1	47.0	42.8
Qwen-7B	73.4	77.9	73.2	71.9	4.6	6.8	55.5	51.4
<b>CoT (0-shot)</b>								
LLaMA-8B	65.3	67.9	63.0	62.5	10.7	8.2	45.8	42.7

Table 5: Macro-F1 across logical relations under 0–3 shot prompting and chain-of-thought (CoT). NV = non-human-validated labels; HV = human-validated subset. NN = NEITHER/NOR, MIX = MIXED.

deeper limitations highlighting how models fail to compose plausibility across logical operators.

### Conjunctive and Disjunctive Reasoning

Across both AND and OR settings, the dominant failure mode is the model’s tendency to recognize only a *single* plausible statement and ignore how plausibility must be composed across clauses. In AND questions, this appears as *single-statement dominance*: the model anchors on a highly plausible clause and fails to reject an implausible partner. For example, in the combustibility question, it predicts *ignites rapidly AND melts into liquid*, correctly identifying ignition while overlooking that combustible materials do not melt.

The same underlying behavior manifests in OR questions. Rather than exploiting the requirement that at least one clause be plausible, the model treats OR as indicating thematic similarity. In the combustibility example, it predicts *melts into liquid OR burns without oxygen*, selecting two implausible but fire-related statements. These errors indicate that models do not enforce the compositional

Op	Failure Pattern	Representative Error (Gold-Pred)
<b>AND</b>	Single-statement dominance	<b>Q:</b> Heat applied to combustible materials? Gold: emits toxic fumes AND decomposes into gases Pred: ignites rapidly AND melts into liquid
<b>OR</b>	Plausibility collapse	<b>Q:</b> Heat applied to combustible materials? Gold: ignites rapidly OR melts into liquid Pred: melts into liquid OR burns without oxygen
<b>NEITHER</b>	Negation inversion	<b>Q:</b> Where do you buy a glass of wine? Gold: NEITHER public park NOR home kitchen Pred: NEITHER licensed restaurant NOR home kitchen
	Plausibility dominance	<b>Q:</b> Where might a child with no home go? Gold: NEITHER school event NOR family gathering Pred: NEITHER youth shelter NOR foster care

Table 6: Representative error patterns for LLaMA-3.1-8B (0-shot) across AND, OR, and NEITHER/NOR.

constraint imposed by the operator; instead, they rely on surface-level plausibility cues and semantic similarity.

### Negation-Based Reasoning

NEITHER/NOR proves most challenging. We observe *negation inversion*, where the model selects the *most plausible* pair of statements despite the operator requiring both to be implausible. For example, in the wine-purchase question, the model predicts *NEITHER licensed restaurant NOR home kitchen*, incorrectly negating the true location. A second pattern, *plausibility dominance*, reflects the model’s tendency to retain highly plausible statements even under negation. For example, in the homeless-child example, the model selects *NEITHER youth shelter NOR foster care*, two locations that are plausible options. These behaviors indicate a broader difficulty with composing negation and plausibility simultaneously.