Towards Consistent Natural-Language Explanations via Explanation-Consistency Finetuning

Anonymous ACL submission

Abstract

Large language models (LLMs) often generate convincing, fluent explanations. However, different from humans, they often generate inconsistent explanations on different inputs. For example, an LLM may generate the explanation "all birds can fly" when answering the question "Can sparrows fly?" but meanwhile answer "no" to the related question "Can penguins fly?". Explanations should be consistent across related examples so that they allow a human to simulate the LLM's decision process on multiple 011 examples. We propose explanation-consistency finetuning (EC-finetuning), a method that adapts LLMs to generate more consistent naturallanguage explanations on related examples. 017

EC-finetuning involves finetuning LLMs on synthetic data that is carefully constructed to contain consistent explanations. Across a variety of question-answering datasets in various domains, EC-finetuning yields a **10.0%** relative explanation consistency improvement on four finetuning datasets, and generalizes to seven out-of-distribution datasets not seen during finetuning (**+4.5%** relative). We will make our code available for reproducibility.

1 Introduction

021

027

037

041

Pre-trained large language models (LLMs) have shown impressive proficiency in a range of complex natural language processing tasks, significantly advancing the field and opening new frontiers for applications (Brown et al., 2020; Touvron et al., 2023; OpenAI, 2023). However, the opaqueness of these models' decision making process has hindered their use in high-stakes applications such as healthcare, and raised issues related to regulatory pressure, safety, and alignment (Goodman and Flaxman, 2016; Amodei et al., 2016; Gabriel, 2020). Moreover, this lack of interpretability has heavily limited the use of LLMs in fields such as social science and data analysis (Ziems et al., 2023),



Figure 1: EC-finetuning adapts an LLM to provide explanations that are more consistent with a user's expectation of LLM answers across related questions.

where trustworthy interpretation (rather than model deployment) is itself the end goal.

One growing avenue into interpretability is natural-language explanations produced by LLMs. These explanations are potentially very powerful, helping users predict model behavior (Johnson-Laird, 1980; Bansal et al., 2019), which is useful to calibrate a model's capacity and limitations, e.g. limiting its demographic bias (Vig et al., 2020). However, these natural-language explanations are limited by the inherent inconsistency of LLMs. For example, one recent work finds that modern LLMs often generate inconsistent explanations and answers on very related questions (Chen et al., 2023b). In fact, LLMs often even struggle to consistently



Figure 2: EC-finetuning synthetically augments the examples in a dataset using LLMs. We instruct the LLM to first generate follow-up questions related to the initial (question, explanation) example, and then to answer the follow-up questions in a manner that is consistent with the explanation of the initial example.

answer rephrasings of the same question (Sclar et al., 2023; Zhang et al., 2023). It is unclear if popular methods for adapting LLMs, e.g. supervised finetuning or reinforcement learning from human feedback, are able to solve this issue.

We address this issue by introducing explanationconsistency finetuning (EC-finetuning). ECfinetuning finetunes an LLM on synthetic data that is precisely constructed to contain consistent explanations. We start with a question-explanation pair (e.g., "*Can sparrows fly?*", "*all birds can fly*"), generate a set of related questions (e.g., "*Can penguins fly?*"), and then answer the related questions to be consistent with the initial explanation (e.g., "*all birds can fly so penguins can fly*"). We generate synthetic data by prompting LLMs, which can be the same as or different from the explanation LLM.

We apply EC-finetuning to question-answering datasets and find that it improves the consistency of natural-language explanations of LLaMA2-13B by **10.0%** relative on four finetuning datasets, and also generalizes to seven out-of-distribution datasets unseen during finetuning (**+4.5%** relative). This suggests that EC-finetuning may be generally useful for helping users build mental models of an LLM from its explanations (see Fig. 1).

2 Related work

061

062

063

065

066

071

079

Generating and improving natural-language explanations A great deal of recent work has studied improving LLM controllability, e.g. through supervised finetuning (Liu et al., 2019), reinforcement learning from human feedback (Ouyang et al., 2022) or learning from human explanations (Stiennon et al., 2020). Two related works study the consistency in the generations made by an LLM, either between the generation and validation of LLMs (Li et al., 2023) or LLM predictions on implications of an original question (Akyürek et al., 2024). In contrast to EC-finetuning, these works do not for

cus directly on improving an LLM's explanation capabilities.

096

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

Many works extend and analyze explanations given by chain-of-thought prompting (Wei et al., 2022), e.g. by evaluating counterfactuals introduced into the chain of thought (Gat et al., 2023), testing their robustness to mistakes introduced into the reasoning chain (Lanham et al., 2023), or using contrastive chain-of-thought to induce reliance on the reasoning chain (Chia et al., 2023). These methods do not alter the underlying LLM, and thus can be used in conjunction with EC-finetuning.

Evaluating natural-language explanations We summarize three existing orthogonal metrics for explanations: consistency, plausibility, and faithfulness. *Consistency*, which we focus on in this work, measures if the model generates consistent explanations on similar examples (Hase and Bansal, 2020; Chen et al., 2023b). *Plausibility* evaluates humans' preference of an explanation based on its factual correctness and logical coherence (Herman, 2017; Lage et al., 2019; Jacovi and Goldberg, 2020). It is different from *faithfulness*, which measures whether an explanation is consistent with the model's internal decision process (Harrington et al., 1985; Jacovi and Goldberg, 2020).

3 Method: EC-finetuning

Explanation-consistency Finetuning ECfinetuning is an intuitive method that augments data in a manner that enhances explanation consistency (Fig. 2). Specifically, it prompts LLMs to augment data in two steps. In the first step, a question-explanation pair is given to an LLM (e.g., "*Can sparrows fly?*", "*all birds can fly*"), with the task of generating follow-up questions related to the explanation of the initial question (e.g., "*Can penguins fly?*"). This is achieved by explicitly prompting the LLM to generate questions that are answerable given the initial explanation. In the second step, another LLM generates answers and explanations for the follow-up questions. To ensure these answers and explanations are consistent with the explanation in the initial question, the initial question-explanation pair is presented in the prompt, alongside explicit instructions to keep the new explanation consistent with the initial (e.g., "*all birds can fly so penguins can fly*".) Finally, these augmented questions, along with their explanations and answers are used for finetuning an LLM to generate more consistent explanations.

135

136

137

138

140

141

142

143

144

145

146

147

148

149

150

152

153

155

156

157

158

159

161

162

163

165

166

168

169

170

171

172

173

174

175

177

178

179

181

182

We use different LLMs for the two data augmentation steps (here, GPT-4 (OpenAI, 2023) for the first step and Claude-2¹ for the second step) to avoid issues with LLMs that favor their own outputs (Zheng et al., 2023). Precise prompts are given in Appendix A.1. Note that the two data augmentation LLMs can be, but are not required to be, identical to the explanation LLM.

Measuring consistency Evaluating the consistency of model explanations is challenging. Here, we follow the metric proposed by Chen et al. 2023b, which measures explanation consistency² as the fraction of answers on follow-up questions that match a human user's expectation based on an explanation (similar to Fig. 1); the metric ranges from 0 to 1, with 1 being a perfect score. Additionally, following Chen et al. 2023b we use LLMs as a user simulator, as it was found to reliably emulate a user's predictions on follow-up questions.

> We evaluate consistency on two types of followup questions: related questions and rephrased questions. For related questions, similar to how we generate EC data, we again use GPT-4 to generate. For rephrased questions, we prompt GPT-4 to generate exact paraphrases of the original questions.

> To ensure this metric from Chen et al. 2023b is reliable, we conduct two sanity checks. First, we check if the metric is stable with respect to how the metric is computed. We run various perturbations (e.g. varying how the followup questions are generated or the explanation format) and find the metric to be quite stable (see Table A1). Second, we study the correlation between explanation consistency and explanation length to see if the metric can be easily hacked by generating shorter/longer explanations. We do not observe a correlation on any of the 7 unseen datasets (see Table A2).

			Consi	Accuracy			
		Relate	ed Qns	Rephra	sed Qns		
		Bl	EC	Bl	EC	Bl	EC
ω	StrategyQA	76.0±1	83.8±1	-	-	73.1±2	74.4±2
nin	MedQA-Diff	$60.8{\pm}1$	$73.1{\pm}1$	_	-	52.9±2	$56.6{\pm}2$
etu	MedQA-Sim	$62.0{\pm}1$	$65.0{\pm}1$	_	-	53.1±2	$55.1{\pm}2$
ine	MedMCQA	$68.2{\pm}1$	$71.3{\pm}1$	_	-	45.5±2	$49.2{\pm}2$
щ	AVG	66.7	73.3	-	-	56.1	58.8
	BoolQ	72.2±1	75.1±1	73.9±1	72.3±1	56.1±2	62.9±2
	NQ	73.1 ± 1	$76.1{\pm}1$	74.0 ± 1	$74.2{\pm}1$	60.8 ± 2	$68.8{\pm}2$
п	MS-Marco	$79.0{\pm}1$	$81.7{\pm}1$	72.2 ± 1	77.6 ± 1	60.0 ± 2	$68.7{\pm}2$
Unsee	OBQA	77.4 ± 1	$80.3{\pm}1$	74.5 ± 1	$74.6{\pm}1$	78.9±2	$74.9{\pm}2$
	MMLU-Med	73.7 ± 1	$77.4{\pm}1$	73.9 ± 1	$80.0{\pm}1$	71.4±2	$68.7{\pm}2$
	PubMedQA	79.5 ± 1	$84.0{\pm}1$	70.2 ± 1	$81.9{\pm}1$	59.3±2	$61.1{\pm}2$
	ARC-Easy	77.5 ± 1	$80.7{\pm}1$	77.1±1	$79.0{\pm}1$	79.2±2	$76.3{\pm}2$
	AVG	76.1	79.4	73.7	77.1	66.5	68.8

Table 1: Main result: EC-finetuning improves explanation consistency (on both related questions and rephrased questions) and prediction accuracy. Error bars show standard error of the mean. Bl: Baseline.

4 Results

4.1 Experimental setup

We perform EC-finetuning on the LLaMA-2 13billion parameter model (Touvron et al., 2023). For finetuning, we use 4 datasets: StrategyQA (Geva et al., 2021), MedMCQA (Pal et al., 2022), and two versions of MedQA (Zhang et al., 2018): MedQA-Sim contains related questions on diagnosis and treatment (similar to the original questions), whereas MedQA-Diff contains related questions on medical facts derived from the original questions.

We additionally evaluate consistency 7 datasets not used for finetuning: on BoolQ (Clark et al., 2019), Natural Questions (NQ) (Kwiatkowski et al., 2019), MS-Marco (Nguyen et al., 2016), OBQA (Mihaylov et al., 2018), MMLU-Medical (Hendrycks et al., 2020), PubMedQA (Jin et al., 2019) and ARC-Easy (Clark et al., 2018). For a cleaner evaluation, these 7 datasets are all converted to have a shared yes-no answer format. We show each dataset's domain and the skills it tests in Table A3. The testing datasets introduce a distribution shift as they cover new domains (science) and new skills (commonsense reasoning and quantitative reasoning) not seen during finetuning. Table A4 shows the size of each dataset.

4.2 Main result: EC-finetuning improves explanation consistency

Table 1 shows the main results for EC-finetuning.EC-finetuning can effectively improve consistency,

185

186

187

188

189

190

191

192

193

194

195

197

198

199

201

202

203

204

205

206

207

210

211

212

213

¹https://www.anthropic.com/index/claude-2

²What we call "consistency", Chen et al. call "counterfactual simulatability precision".

	BQ	PMQA	NQ	MM	ARC	OBQA	MMLU	AVG
Baseline	0.05	-0.05	0.14	0.13	0.15	0.22	0.05	0.10
EC	0.23	0.04	0.30	0.27	0.13	0.26	0.07	0.19

Table 2: Pearson Correlation between explanation consistency and prediction accuracy across examples. Inconsistent explanations suggest wrong predictions.

	BQ	PMQA	NQ	MM	ARC	OBQA	MMLU	AVG
Correct	5.2%	8.3%	5.9%	4.9%	4.8%	5.5%	4.8%	5.8%
Incorrect	0.0%	-5.6%	-2.8%	1.6%	6.5%	5.1%	3.6%	0.8%

Table 3: Consistency improvement from EC on correct vs. incorrect predictions. EC-finetuning improves explanation consistency more on correct predictions.

yielding an average relative improvement of **10.0%** 214 for tasks seen during finetuning and 4.5% for un-215 seen tasks. An improvement is seen for every 216 dataset studied here and for both types of followup 217 questions. The largest gain in consistency after ECfinetuning is for MedQA-Diff; this suggests that 219 EC-finetuning can also improve the LLM's explanation consistency on related questions that are more 221 different from the original questions. These consistency improvements also come with modest accu-223 racy improvements (5.2% relative for finetuning tasks and 4.3% relative for unseen tasks). There 225 is no significant correlation between improvement in consistency and the improvement in accuracy 227 (Pearson correlation coefficient $\rho = 0.001$). This 228 suggests that the consistency improvement derived from EC-finetuning differs from the improvement attained by standard supervised finetuning.

EC-finetuning using only a single LLM We explore a simplified setting, where EC-finetuning is 233 run using only the LLaMA-2 13-billion parameter, both for synthetic data generation and explanation finetuning. This setting tests whether ECfinetuning can be used with smaller LLMs and 237 whether those LLMs can improve their own expla-238 nation consistency. We find that when running ECfinetuning on StrategyQA, EC-finetuning yields a 4.4% relative improvement but decreases accuracy 241 by 5.4%. This suggests that EC-finetuning may succeed in improving explanation consistency in 243 today's relatively small models, but can incur some tradeoffs as a result, i.e. decreasing accuracy. 245

4.3 Analysis

246

EC-finetuning improves explanation consistency
in different ways. Table A5 shows examples of

explanations before/after EC-finetuning. The consistency of the explanation in both examples increases after EC-finetuning, but in different ways. In the first example, EC-finetuning encourages the model to generate more precise explanations that are not overgeneralized/vague. On the other hand, in the second example, EC-finetuning does not change the explanation the model generates for the initial question, but instead changes the model's predictions on related questions to be more consistent with the explanation on the initial question. 249

250

251

252

253

254

255

256

257

258

259

260

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

285

286

289

290

291

292

293

294

295

296

297

Inconsistent explanations suggest incorrect predictions. Do LLMs generate more consistent explanations on correct predictions? We study the correlation between explanation consistency and prediction accuracy across different examples of the same dataset. We find that the baseline model shows a positive correlation of 0.099 (Pearson), and this correlation increases to 0.185 after ECfinetuning (dataset-level breakdown in Table 2). This indicates that inconsistent explanations suggest wrong predictions, and we may calibrate LM's predictions based on the consistency of its explanations (Chen et al., 2023a).

EC-finetuning improves consistency more on correct predictions. We compare the consistency improvement from EC-finetuning on correct versus incorrect predictions. EC-finetuning improves explanation consistency on correct predictions by 5.7% relative but only 1.2% relative on incorrect predictions (see full breakdown in Table 3). This is possibly because examples where the model makes mistakes are likely challenging ones beyond the model's capability, so the model can produce inconsistent hallucinations.

5 Conclusion

Forming reliable mental models of LLMs will be critical to understanding and deploying them as their capabilities continue to grow. EC-finetuning is a step in this direction, helping to align models to not only be more accurate, but also to provide more self-consistent explanations. A big part of future work will be *i*) scaling up EC-finetuning to larger models, *ii*) exploring EC-finetuning on more complicated tasks (e.g., multimodal understanding), *iii*) studying how LLMs represent consistency in their parameters and if EC-finetuning improves this representation, and *iv*) proposing other methods to improve explanation consistency.

6 Limitations

298

301

302

304

306

307

311

312

313

314

316

319

320

321

322

325

326

328

330

333

335

338

340

341

342

345

346

347

351

352

Our work has several limitations that we leave to future work. First, due to limited computational resources we only experiment with models up to 13 billion parameters. Future work should scale up our experiments to larger LMs. Second, our work focuses on explanations in the text modality. Future work should extend EC-finetuning to multimodal explanations.

References

- Afra Feyza Akyürek, Ekin Akyürek, Leshem Choshen, Derry Wijaya, and Jacob Andreas. 2024. Deductive closure training of language models for coherence, accuracy, and updatability. arXiv preprint arXiv:2401.08574.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. Concrete problems in ai safety. arXiv preprint arXiv:1606.06565.
- Gagan Bansal, Besmira Nushi, Ece Kamar, Walter S Lasecki, Daniel S Weld, and Eric Horvitz. 2019. Beyond accuracy: The role of mental models in human-ai team performance. In Proceedings of the AAAI conference on human computation and crowdsourcing.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901.
- Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, and He He. 2023a. On the relation between sensitivity and accuracy in in-context learning. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 155-167, Singapore. Association for Computational Linguistics.
- Yanda Chen, Ruiqi Zhong, Narutatsu Ri, Chen Zhao, He He, Jacob Steinhardt, Zhou Yu, and Kathleen McKeown. 2023b. Do models explain themselves? counterfactual simulatability of natural language explanations. arXiv preprint arXiv:2307.08678.
- Yew Ken Chia, Guizhen Chen, Luu Anh Tuan, Soujanya Poria, and Lidong Bing. 2023. Contrastive chain-of-thought prompting.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. arXiv preprint arXiv:1905.10044.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. arXiv preprint arXiv:1803.05457.
- Iason Gabriel. 2020. Artificial intelligence, values, and alignment. Minds and machines, 30(3):411-437.
- Yair Gat, Nitay Calderon, Amir Feder, Alexander Chapanin, Amit Sharma, and Roi Reichart. 2023. Faithful explanations of black-box nlp models using llm-generated counterfactuals. arXiv preprint arXiv:2310.00603.

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan	353
Roth, and Jonathan Berant. 2021. Did aristotle use a lap-	354
top? a question answering benchmark with implicit reason-	355
ing strategies. <i>Transactions of the Association for Compu-</i>	356
<i>tational Linguistics</i> .	357
Bryce Goodman and Seth Flaxman. 2016. European union regulations on algorithmic decision-making and a" right to explanation". <i>arXiv preprint arXiv:1606.08813</i> .	358 359 360
Leo A Harrington, Michael D Morley, A Šcedrov, and	361
Stephen G Simpson. 1985. <i>Harvey Friedman's research</i>	362
<i>on the foundations of mathematics.</i>	363
Peter Hase and Mohit Bansal. 2020. Evaluating explain-	364
able AI: Which algorithmic explanations help users predict	365
model behavior? In <i>Proceedings of the Association for</i>	366
<i>Computational Linguistics</i> .	367
Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Man-	368
tas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Mea-	369
suring massive multitask language understanding. <i>arXiv</i>	370
<i>preprint arXiv:2009.03300.</i>	371
Bernease Herman. 2017. The promise and peril of human evaluation for model interpretability. <i>ArXiv</i> .	372 373
Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully	374
interpretable NLP systems: How should we define and	375
evaluate faithfulness? In <i>Proceedings of the Association</i>	376
<i>for Computational Linguistics</i> .	377
Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. <i>arXiv preprint arXiv:1909.06146</i> .	378 379 380 381
Philip N Johnson-Laird. 1980. Mental models in cognitive science. <i>Cognitive science</i> .	382 383
Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield,	384
Michael Collins, Ankur Parikh, Chris Alberti, Danielle	385
Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al.	386
2019. Natural questions: a benchmark for question answer-	387
ing research. <i>Transactions of the Association for Computa-</i>	388
<i>tional Linguistics</i> , 7:453–466.	389
Isaac Lage, Emily Chen, Jeffrey He, Menaka Narayanan, Been Kim, Sam Gershman, and Finale Doshi-Velez. 2019. An evaluation of the human-interpretability of explanation. <i>ArXiv</i> .	390 391 392 393
Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit	394
Steiner, Carson Denison, Danny Hernandez, Dustin Li,	395
Esin Durmus, Evan Hubinger, Jackson Kernion, et al.	396
2023. Measuring faithfulness in chain-of-thought reason-	397
ing. <i>arXiv preprint arXiv:2307.13702</i> .	398
Xiang Lisa Li, Vaishnavi Shrivastava, Siyan Li, Tatsunori	399
Hashimoto, and Percy Liang. 2023. Benchmarking and im-	400
proving generator-validator consistency of language mod-	401
els. <i>arXiv preprint arXiv:2310.01846</i> .	402
Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019. Multi-task deep neural networks for natural language understanding. <i>arXiv preprint arXiv:1901.11504</i> .	403 404 405
Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sab-	406
harwal. 2018. Can a suit of armor conduct electricity? a	407
new dataset for open book question answering. In <i>EMNLP</i> .	408

Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. Ms marco: A human generated machine reading comprehension dataset. *choice*, 2640:660.

OpenAI. 2023. Gpt-4 technical report. ArXiv.

409

410

411

412

413 414

415

416

417

418

419

420 421

422

423

494

425

426

427

428

429

430

431

432

433

434

435 436

437

438

439

440

441

442

443

444

445

446

447

448

449 450

451

452

453

454

455

456

457 458

459

460

461

462

463

464

465

466

467

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. Advances in Neural Information Processing Systems.
- Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multisubject multi-choice dataset for medical domain question answering. In *Conference on Health, Inference, and Learning*, pages 248–260. PMLR.
 - Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*.
 - Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021.
 - Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
 - Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. In Advances in Neural Information Processing Systems.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
 - Qingru Zhang, Chandan Singh, Liyuan Liu, Xiaodong Liu, Bin Yu, Jianfeng Gao, and Tuo Zhao. 2023. Tell your model where to attend: Post-hoc attention steering for llms. *arXiv preprint arXiv:2311.02262*.
 - Xiao Zhang, Ji Wu, Zhiyang He, Xien Liu, and Ying Su. 2018. Medical exam question answering with large-scale reading comprehension. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
 - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MTbench and chatbot arena. In *Thirty-seventh Conference* on Neural Information Processing Systems Datasets and Benchmarks Track.
- Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2023. Can large language models transform computational social science? *arXiv preprint arXiv:2305.03514*.

568

569

518

A Appendix

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

496

497

498

499

502

503

505

506

510

511

512

513

514

515

516

517

A.1 Prompting details

We prompt LLMs to generate the EC training data in two steps 1) generate follow-up questions related to an initial explanation, 2) answer the related questions conditioned on the initial explanation. We include a few demonstration examples for both steps (in-context learning), which are sampled from each dataset. We include below the prompts we use for the StrategyQA dataset, and will make the demonstration examples for other datasets publicly available.

Follow-up Questions Generation Human: In the questions below, you will be asked to read a starter question and its answer. After that you will be asked to write a follow-up question that can be answered based on the starter QA, and write your answer to the follow-up question based on the starter QA. Your follow-up question should be self-contained even without the starter question.

Assistant: here is my response. okay.

Human: Starter Question: Can the Palace of Westminster tell time in the dark? Answer to the Starter Question: The Palace of Westminster is a building. Buildings cannot tell time.

Assistant: here is my response. Follow-up Question: Can the White House tell time? Answer to the Follow-up Question: We know from the starter QA that buildings cannot tell time. Since the White House is also a building, the White House cannot tell time either. So the answer is no.

Human: Starter Question: Are psychiatric patients welcome to join the United States Air Force? Answer to the Starter Question: The United States Air Force is a branch of the United States military. The United States military does not allow psychiatric patients to join.

Assistant: here is my response. Follow-up Question: Are psychiatric patients welcome to join the United States Army? Answer to the Follow-up Question: We know from the starter QA that the United States military does not allow psychiatric patients to join. Because United States Army is a branch of the United States military, psychiatric patients are not welcome to join the United States Army either. So the answer is no.

Human: Starter Question: Does The Hague border multiple bodies of water? Answer to the Starter Question: The Hague is a city in the Netherlands. The Netherlands borders the North Sea. Assistant: here is my response. Follow-up Question: Is Amsterdam near the North Sea? Answer to the Follow-up Question: We know from the starter QA that the Netherlands borders the North Sea. Because Amsterdam is a city in the Netherlands, Amsterdam is also near the North Sea. So the answer is yes.

Human: Starter Question: Can an Asian black bear use chopsticks? Answer to the Starter Question: Asian black bears are omnivores. Omnivores can use chopsticks.

Assistant: here is my response. Follow-up Question: Can pigs use chopsticks? Answer to the Follow-up Question: We know from the starter QA that omnivores can use chopsticks. Because pigs are omnivores, pigs can use chopsticks as well. So the answer is yes.

Human: Starter Question: Does open heart surgery finish before entirety of American Ballet Theatre's Swan Lake? Answer to the Starter Question: The American Ballet theatre's Swan Lake has a run time of 45 minutes. The National Heart, Lung, and Blood Institute states that a coronary artery bypass takes 3 to 6 hours.

Assistant: here is my response. Follow-up Question: Is it possible to watch the entirety of American Ballet Theatre's Swan Lake 2 times before an open heart surgery finishes? Answer to the Followup Question: We know from the starter QA that the run time of Swan Lake is 45 minutes and an open heart surgery takes 3 to 6 hours. Therefore it is possible to watch Swan Lake twice before the surgery finishes (90 minutes). So the answer is yes.

Human: Starter Question: Could a bee hummingbird balance a scale with a single pea on it? Answer to the Starter Question: The bee hummingbird is the smallest bird in the world. The bee hummingbird weighs about 2 grams. A pea weighs about 1 gram.

Assistant: here is my response. Follow-up Question: Could a pea balance a scale with a dollar bill on it? Answer to the Follow-up Question: We know from the starter QA that a pea weighs about 1 gram. Because a dollar bill weighs about 1 gram as well, a pea can balance a scale with a dollar bill on it. So the answer is yes.

Human: Starter Question: While viewing "Scary Movie" is the viewer likely to experience an increase in adrenaline? Answer to the Starter Question: The movie "Scary Movie" is a comedy. Comedies are not scary.

Assistant: here is my response. Follow-up Ques-

Explanation Format		Followup Question Generation			Consistency
Brevity	Includes Answer	Model	Decoding	Prompt	
Facts	Yes	GPT-4	List	А	84.8 ± 2
Facts	Yes	GPT-4	List	В	$88.8{\pm}2$
Facts	Yes	GPT-4	Sampling	А	85.1 ± 2
Facts	Yes	GPT-4	Sampling	В	87.6 ± 2
Facts	No	GPT-4	List	А	$85.7{\pm}2$
Facts	No	GPT-4	List	В	$85.9{\pm}2$
Facts	No	GPT-4	Sampling	А	$82.9{\pm}3$
Facts	No	GPT-4	Sampling	В	$84.9{\pm}3$
Facts+Reasoning	Yes	GPT-4	List	А	85.6 ± 2
Facts+Reasoning	Yes	GPT-4	List	В	85.1 ± 2
Facts+Reasoning	Yes	GPT-4	Sampling	А	$87.7{\pm}2$
Facts+Reasoning	Yes	GPT-4	Sampling	В	$84.9{\pm}3$
Facts+Reasoning	No	GPT-4	List	А	$85.5{\pm}2$
Facts+Reasoning	Yes	GPT-3	Sampling	А	$83.2{\pm}2$

Table A1: We evaluate the stability of the consistency metric with respect to how the metric is computed. We run several perturbations for the explanation format (whether the explanation includes the reasoning chain besides listing relevant facts, and whether the explanation includes the answer) and how the followup questions are generated (model, decoding, and prompt variations). We find that the consistency metric is quite stable with respect to how the metric is computed.

	Baseline	EC
BoolQ	-0.059	0.039
PubMedQA	0.044	0.007
NQ	-0.002	0.028
MS-Marco	0.084	0.006
ARC-Easy	0.030	0.012
OBQA	-0.066	-0.034
MMLU-Med	-0.020	0.003

Table A2: We find no significant correlation (p-value < 0.005) between explanation consistency and explanation length for either the baseline or EC.

tion: Is "Superbad" scary? Answer to the Followup Question: We know from the starter QA that comedies are not scary. "Superbad" is also a comedy, so "Superbad" is not scary either. So the answer is no.

571

572

573

575

576

578

580

582

583

585

586

Human: Starter Question: Is Mixed martial arts totally original from Roman Colosseum games? Answer to the Starter Question: Mixed Martial arts in the UFC takes place in an enclosed structure called The Octagon. The Roman Colosseum games were fought in enclosed arenas where combatants would fight until the last man was standing. Mixed martial arts contests are stopped when one of the combatants is incapacitated. The Roman Colosseum was performed in front of crowds that numbered in the tens of thousands. Over 56,000 people attended UFC 193.

_			
	Dataset	Domain	Skills
FT	StrategyQA MedQA MedMCQA	Open Medical Medical	Facts, Reason Facts, Reason Facts
Unseen	BoolQ NQ MS-Marco OBQA MMLU-Med PubMedQA ARC-Easy	Open Open Open Medical Medical Science	Facts Facts Facts Reason, <i>Commonsense</i> Facts Facts, Reason, <i>Quantitative</i> Facts, Reason

Table A3: We show the domain and required skills of each dataset. New domains and skills in testing datasets unseen during training are bold and italicized.

Write 10 follow-up yes or no questions. Start each question by "1. ", "2. ", etc. Skip writing answers to your follow-up questions. Remember, your task is to write follow-up questions that can be answered based on the starter QA. Your follow-up questions should be self-contained even without the starter question. 587

588

589

590

591

592

593

594

595

596

597

598

599

600

Follow-up Questions Answering Human: In the questions below, you will be asked to read a starter question and its answer. After that you will be asked to read a follow-up question and judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.".

8

		Train	Dev	Test
FT	StrategyQA MedQA MedMCQA	1400 2498 2500	100 100 100	790 600 600
Unseen	BoolQ NQ MS-Marco OBQA MMLU-Medical PubMedQA ARC-Easy	 	- - - - -	500 500 500 500 500 445 500

Table A4: Sizes of all datasets used.

If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong.

601

606

610

611

613

614

615

616

617

618

619

626

634

635

638

Assistant: here is my response. okay.

Human: Starter Question: Can the Palace of Westminster tell time in the dark? Answer to the Starter Question: The Palace of Westminster is a building. Buildings cannot tell time. Follow-up Question: Can the White House tell time? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question:

Assistant: here is my response. The White House is a building. Buildings cannot tell time. So the answer is no.

Human: Starter Question: Are psychiatric patients welcome to join the United States Air Force? Answer to the Starter Question: The United States Air Force is a branch of the United States military. The United States military does not allow psychiatric patients to join. Follow-up Question: Are psychiatric patients welcome to join the United States Army? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question: 639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

Assistant: here is my response. The United States military does not allow psychiatric patients to join. The United States Army is a branch of the United States military. So the answer is no.

Human: Starter Question: Does The Hague border multiple bodies of water? Answer to the Starter Question: The Hague is a city in the Netherlands. The Netherlands borders the North Sea. Follow-up Question: Is Amsterdam near the North Sea? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question:

Assistant: here is my response. Amsterdam is a city in the Netherlands. Netherlands borders the North Sea. So the answer is yes.

Human: Starter Question: Does Marco Rubio have a close relationship with Allah? Answer to the Starter Question: Marco Rubio is a Christian. Christians do not have a close relationship with Allah. Follow-up Question: Is Barack Obama a Christian? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "Ac-

Initial Question	Baseline Explanation	EC Explanation
Are they making a season 4 of the flash?	The flash is a TV show that has been on for 3 seasons. TV shows are typically cancelled after 3 seasons. Con- sistency: 0.00	The CW has renewed The Flash for a seventh season. Consistency: 0.74
Does age of ultron come after winter soldier?	Winter soldier was released in 2014. Age of Ultron was released in 2015. Consistency: 0.06	Captain America: The Winter Soldier was released in 2014. Avengers: Age of Ultron was released in 2015. Consistency: 0.95

Table A5: Examples of more consistent explanations from EC-finetuning. EC-finetuning improves explanation consistency in different ways.

cording to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question:

693

694

695

696

701

702

703

704

708

710

711

712

713

714

715

717

718

719

720

721

723

725

727

731

Assistant: here is my response. The starter QA does not include information that directly answer this follow-up question. So the answer is unknown.

Human: Starter Question: Can an Asian black bear use chopsticks? Answer to the Starter Question: Asian black bears are omnivores. Omnivores can use chopsticks. Follow-up Question: Can pigs use chopsticks? Judge whether the starter QA directly helps choosing a single answer for the followup question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question:

Assistant: here is my response. Pigs are omnivores. Omnivores can use chopsticks. So the answer is yes.

Human: Starter Question: Does open heart surgery finish before entirety of American Ballet Theatre's Swan Lake? Answer to the Starter Question: The American Ballet theatre's Swan Lake has a run time of 45 minutes. The National Heart, Lung, and Blood Institute states that a coronary artery bypass takes 3 to 6 hours. Follow-up Question: Is it possible to watch the entirety of American Ballet Theatre's Swan Lake 2 times before an open heart surgery finishes? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question:

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

Assistant: here is my response. The run time of Swan Lake is 45 minutes. An open heart surgery takes 3 to 6 hours. 45 minutes * 2 = 90 minutes < 3 hours. So the answer is yes.

Human: Starter Question: Could a bee hummingbird balance a scale with a single pea on it? Answer to the Starter Question: The bee hummingbird is the smallest bird in the world. The bee hummingbird weighs about 2 grams. A pea weighs about 1 gram. Follow-up Question: Could a pea balance a scale with a dollar bill on it? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question:

Assistant: here is my response. A pea weighs about 1 gram. A dollar bill weighs about 1 gram. So the answer is yes.

Human: Starter Question: While viewing "Scary Movie" is the viewer likely to experience an increase in adrenaline? Answer to the Starter Ques-

825

tion: The movie "Scary Movie" is a comedy. Come-773 dies are not scary. Follow-up Question: Is "Super-774 bad" scary? Judge whether the starter QA directly 775 helps choosing a single answer for the follow-up 776 question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to 778 answer the follow-up question, explain your rea-779 soning as clearly and as detailed as possible using all relevant information in the starter QA, end your 781 answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up **Ouestion**:

> Assistant: here is my response. The movie "Superbad" is a comedy. Comedies are not scary. So the answer is no.

790

791

796

797

803

804

807

811

812

813

814

817

818

819 820

821

822

824

Human: Starter Question: Do gorillas fight with panda bears? Answer to the Starter Question: Gorillas are herbivores. Pandas are herbivores. Herbivores do not fight with other herbivores. Followup Question: Do herbivores fight with carnivores? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question:

Assistant: here is my response. The starter QA does not include information that directly answer this follow-up question. So the answer is unknown.

Human: Starter Question: Can citrus grow in Ulaanbaatar? Answer to the Starter Question: Citrus trees can grow in Ulaanbaatar. Thus, citrus can grow in Ulaanbaatar. Follow-up Question: Can 816 palm trees grow in Ulaanbaatar? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer

is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question:

Assistant: here is my response. The starter OA does not include information that directly answer this follow-up question. So the answer is unknown.

Human: Starter Question: Is Mixed martial arts totally original from Roman Colosseum games? Answer to the Starter Question: Mixed Martial arts in the UFC takes place in an enclosed structure called The Octagon. The Roman Colosseum games were fought in enclosed arenas where combatants would fight until the last man was standing. Mixed martial arts contests are stopped when one of the combatants is incapacitated. The Roman Colosseum was performed in front of crowds that numbered in the tens of thousands. Over 56,000 people attended UFC 193. Follow-up Question: Is the UFC Octagon considerably smaller than the Roman Colosseum? Judge whether the starter QA directly helps choosing a single answer for the follow-up question. If not, end your answer with "So the answer is unknown.". If yes, use the starter QA to answer the follow-up question, explain your reasoning as clearly and as detailed as possible using all relevant information in the starter QA, end your answer with "So the answer is yes/no.", and do NOT explicitly mention "the starter QA" or "According to the starter QA" in your answer. Stick to the starter QA when you answer the follow-up question, even if the reasoning or claims in the starter QA are wrong. Answer to the Follow-up Question: