Characterizing Large Language Models as Rationalizers of Knowledge-intensive Tasks

Anonymous ACL submission

Abstract

Large language models (LLMs) are proficient 002 at generating fluent text with minimal taskspecific supervision. However, their ability to generate rationales for knowledge-intensive tasks (KITs) remains under-explored. Generating rationales for KIT solutions, such as commonsense multiple-choice QA, requires ex-007 ternal knowledge to support predictions and refute alternate options. In this work, we consider the task of generating retrieval-augmented rationalization of KIT model predictions via external knowledge-guidance within a fewshot setting. Surprisingly, crowd-workers pre-013 ferred LLM-generated rationales over exist-015 ing crowd-sourced rationales, generated in a similar knowledge-guided setting, on aspects such as factuality, sufficiency, and convincing-017 ness. However, follow-up fine-grained evaluation of such rationales highlight the need for further improvements in conciseness, novelty, and domain invariance. Additionally, through an expert-sourced study evaluating the reliability of the rationales, we demonstrate that humans' trust in LLM-generated rationales erode when communicated faithfully, *i.e.*, without taking model prediction accuracy into account. We 027 find that even instrumenting simple guardrails can be an effective for reliable rationalization.

1 Introduction

037

041

In recent years, generating *rationales* (*i.e.*, freetext explanations) of natural language understanding tasks has been increasingly explored in the field of explainable NLP. Such rationales — while less functionally grounded, *i.e.*, they may not entirely reflect the model's behavior — provide an effective interface to interpretably communicate model decisions to end-users (Hendricks et al., 2016; Camburu et al., 2018; Madsen et al., 2022; Gurrapu et al., 2023). Generating these rationales via direct supervision (Ehsan et al., 2018; Narang et al., 2020) or fine-tuning (Aggarwal et al.,



<corroborate> Commonsense suggests that a waiter, who is generally located in a restaurant, typically presents a bill. Therefore, the answer is "present bill" because this is a common practice at the end of a meal in a restaurant.

<refute> While a waiter can serve food, set the table, and serve a meal, these actions typically occur before or during the meal, not at the end. The option 'eat' is not suitable as it is not a typical duty of a waiter during their service.



Figure 1: a) A commonsense question with multiple choices and knowledge extracted from ConceptNet and b) proposed LLM-generated rationale corroborating the selected answer and refuting the other choices.

2021; Rei et al., 2022) requires the collection of high-quality human-authored rationales. Collecting such rationales via crowd-sourcing is expensive, difficult to standardize, and lacks generalizability to different domains (Wiegreffe and Marasović, 2021; Tan, 2021). Recent work (Wiegreffe et al., 2022) showcases that large language model (LLM) generated rationales, obtained via few-shot in-context learning (Radford et al., 2019; Brown et al., 2020; Huang et al., 2023), alleviate these challenges while showcasing surprising effectiveness over crowdsourced rationales on dimensions such as human preference. However, characterizing the suitability of LLMs as rationalizers of knowledge-intensive task (KIT) decisions such as commonsense question answering (CSQA (Talmor et al., 2019)) and open book question answering (OBQA (Mihaylov et al., 2018)) requires further investigation due to the difference in scope and setting from prior work (Wiegreffe et al., 2022).

Firstly, KITs such as CSQA and OBQA are framed as multiple-choice questions, requiring models to select one answer from several choices (see Figure 1a). Therefore, a corresponding wellformed rationale is required to be (a) comprehensive, *i.e.*, state facts that are not present in the question but are essential for rationalization, and (b) refutation complete, *i.e.*, rationalize why the rest of the choices are incorrect or not best suited as

070

043

the answer (Aggarwal et al., 2021). We show an 071 example of such a rationale in Figure 1b. However, 072 LLM-generated rationales in prior work (Wiegreffe 073 et al., 2022) have only been evaluated on their corroboration capabilities. Secondly, LLM-generated rationales in prior work are abstractive (Gurrapu et al., 2023), lacking grounding on external knowl-077 edge sources crucial for accomplishing the task — KIT models designed for CSQA and OBQA (Feng et al., 2020; Yasunaga et al., 2021, 2022) refer to external sources such as ConceptNet (Speer et al., 2017) (see Figure 1a). Finally, KIT models may predict incorrectly — faithfully rationalizing such mistakes may erode the end-user's trust in the gen-084 erated rationales. Existing approaches in explainable NLP omit the incorrect prediction confounder and evaluate only rationales of correct predictions. However, with LLM-generated rationales being increasingly adopted in real-world scenarios, such as rationalizing why a candidate is suitable for an advertised job¹, it is important to scrutinize the practical implications of such deployments and inform guidelines for safe and responsible adoption.

> Given the setting of generating corroborating and refutation complete rationales of KIT model decisions, we explore the suitability of retrievalaugmented rationale generation using LLMs. We enrich the prompt to LLMs with relevant knowledge retrieved from external sources to condition the rationale generation on facts. More specifically, we generate knowledge-guided rationales contain-similar to Figure 1b — via few-shot prompting of LLMs. We conducted three human subject studies to evaluate the effectiveness of such rationales in communicating KIT model decisions. The observations from these studies enable coarse- and fine-grained characterization of the strengths and weaknesses of LLM-generated knowledge-guided rationalization of KIT model decisions.

100

101

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

More specificlly, we conduct two studies via crowdsourcing to evaluate the preferability and acceptability of such rationales to crowd-workers. In another study involving experts — motivated by existing literature on trust in explainable AI (Hoffman et al., 2018; Stites et al., 2021) — we explore the implications of faithfully rationalizing KIT model decisions irrespective of their correctness. The crowd-sourced studies demonstrate that, more often than not, crowdworkers prefer LLM-generated rationales to crowdsourced rationales in existing datasets, citing their factuality, sufficiency, and convincing refutation. Follow-up fine-grained analysis reveals that LLM-generated rationales still have significant room for improvement along dimensions such as insightfulness (i.e., providing new information), redundancy (i.e., avoiding repetitive text), and generalizability (i.e., domain invariance.) The expert-sourced study confirms that faithful rationalization of incorrect model predictions degrades humans' trust in the generated rationales. We further explore the utility of instrumenting mechanisms to intervene the incorrect predictions via a reviewthen-rationalize pipeline instead of faithfully rationalizing and find that even simple strategies may help intervene up to 71% of the incorrect predictions. We will publicly release the code and data.

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

2 Knowledge-enhanced Rationalization

KIT models such as MHGRN (Feng et al., 2020), QAGNN (Yasunaga et al., 2021), and Dragon (Yasunaga et al., 2022) combine language model and knowledge graph representations to solve complex tasks such as commonsense QA (Talmor et al., 2019). We aim to generate rationales that corroborate the KIT model's prediction with additional relevant facts while refuting the other choices (see Figure 1.) Our approach is similar to existing retrievalaugmented generation strategies with LLMs (Peng et al., 2023; Lazaridou et al., 2022; Zhao et al., 2023; Mei et al., 2023). To guide the generation of these rationale components, i.e., corroboration and refutation, we retrieve facts concerning the knowledge-intensive task — e.g., questions and choices in CSQA and OBQA - from a knowledge graph such as ConceptNet (Speer et al., 2017). We then prompt an LLM to rationalize the prediction via conditioning on the provided knowledge. Figure 2 outlines the rationalization process given an input, i.e., question, choices, and model prediction.

Given an external knowledge-graph such as ConceptNet (Speer et al., 2017), we employ the knowledge extraction strategy used in QAGNN (Yasunaga et al., 2021) to first retrieve the facts relevant to a question and then select top-k (k = 5) facts based on their RoBERTa (Liu et al., 2019) score given the question and a choice. Such selection enables us to fit the knowledge facts within the token limits of an LLM prompt. We employ greedy decoding-based few-shot prompting to query an

¹https://www.businessinsider.com/sc/

indeed-is-embracing-ai-to-power-the-future-of-work



Figure 2: Given an Input (*i.e.*, QA and model prediction), an LLM is prompted to generate a rationale with few-shot examples sampled from an expert-written pool.

LLM for rationalization. Each example in the prompt contains a QA task, the corresponding KIT model prediction, facts retrieved from Concept-Net, and expert written rationale corrborating the prediction and refuting other choices. We opted for expert-authored rationales due to their reported effectiveness over crowdsourced rationales (Wiegreffe et al., 2022). The paper's authors collaboratively crafted high-quality rationales to compile the expert-written pool. We provide a detailed description of the prompt design in Appendix A (Table 4.) Give a new multiple-choice question; we combine the question, the model prediction, and the corresponding extracted facts with the few-shot examples sampled from the expert-pool to formulate the final prompt (see Figure 2.)

170

171

172

173

174

175

176

177

178

179

180

181

183

184

186

187

191

192

193

194

195

197

198

199

3 Evaluation of Rationales

Due to a lack of suitable automated methods for evaluating the rationale quality (Clinciu et al., 2021; Kayser et al., 2021) and credibility, we conducted three studies to address the following questions:

RQ1. How effective are the LLM-generated rationales in communicating KIT model decisions compared to crowdsourced rationales? (§ 4)

RQ2. To what degree do the fine-grained rationale characteristics influence its effectiveness and how generalizable are these observations? (\S 5)

RQ3. How does faithful rationalization of model predictions impact humans' trust in the LLM-generated rationales? (§ 6)

200Datasets and Prompts. We select QAGNN (Ya-201sunaga et al., 2021) as the KIT model due to its202well-documented code repository and availability203of pre-trained model weights. We consider two204datasets of multiple-choice QA tasks related to205commonsense knowledge, CSQA (Talmor et al.,2062019), and elementary-level science, OBQA (Mi-207haylov et al., 2018). Following the existing KIT208models, we use ConceptNet (Speer et al., 2017) as

our external knowledge source. For both datasets, we report results on a fixed, randomly-sampled 250instance test set. We sample these instances from the test set prepared for these datasets (Feng et al., 2020). We employed GPT-3.5 text-davinci-003 (temperature = 0) as the LLM rationalizer. We randomly selected 40 instances from each of the CSQA and OBQA datasets — different from the 250 test instances — to be included in the expertwritten example pool. See Appendix A for details.

209

210

211

212

213

214

215

216

217

218

219

220

221

222

224

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

Faithful Rationalization Studies. We conducted two crowdsourced studies aimed at addressing RQ1 and RQ2. For both studies, we only consider rationalization of correct KIT model predictions, i.e., faithful rationalization. The approach is similar to prior work (Aggarwal et al., 2021; Wiegreffe et al., 2022; Marasovic et al., 2022; Kayser et al., 2021) that also removed the confounder, *i.e.*, rationalization of incorrect model prediction, by only considering rationales for correctly predicted instances. We used Amazon Mechanical Turk for crowdsourcing evaluation. For HITs in both studies, we asked targeted questions to obtain coarseand fine-grained feedback on the rationales of a KIT model decision. We detail these evaluation metrics in the respective sections discussing the studies. Due to the subjectivity of some of the instances of the CSQA dataset, following Wiegreffe (Wiegreffe et al., 2022), we instruct workers for both the studies to consider the KIT model prediction to be correct even if they disagree with it. We undertook several quality control measures from vetting and recruitment of crowdworkers to accounting for order effect of tasks and individual annotator bias. Besides detailing these measures, we include the study interface design and additional statistical information in Appendices B and E.

Credible Rationalization Study. To address *RQ3*, inspired by existing work on trust in explainable AI (Hoffman et al., 2018; Stites et al., 2021; Smith-Renner et al., 2020), we conducted a confirmatory study in the context of explainable NLP (*i.e.*, LLM-generated rationalization) to explore credibility of rationales on aspects such as agreement, confidence, reliability, and user satisfaction, among others. In this study, we consider rationales generated on both correct and incorrect KIT model predictions. The study was conducted via a Slack campaign within Company X, an industrial research lab, with NLP, data management, and machine learning as the primary research areas.

4 LLMs vs Humans as Rationalizers

260

262

263

265 266

267

273

274

275

276

279

287

289

290

294

301

309

We first compare LLM-generated rationales of the CSQA (Talmor et al., 2019) tasks with corresponding crowdsourced rationales from ECQA dataset (Aggarwal et al., 2021). The ECQA rationales are similar in construct to our setting containing corroboration and refutation of CSQA tasks. We exclude CoS-E (Rajani et al., 2019), another crowdsourced free-text rationales dataset, as those rationales are not refutation complete. Moreover, ECQA rationales are reported to be overall better than CoS-E in rationalizing KIT decisions (Aggarwal et al., 2021; Sun et al., 2022). We explain the dataset selection criteria in further detail in Appendix A. Following are the key study takeaways:

- knowledge-guided rationales are preferable (67.2%) to crowdworkers compared to crowdwritten rationales, while showcasing substantial increase in preference (45.7%) than prior work (Wiegreffe et al., 2022).
- fine-grained aspects of a rationale such as supportiveness, sufficiency, and convincingness weakly predict such preferences.

4.1 Study Setting

In each of the 250 HITs (three different crowdworkers per HIT), a crowd-worker was presented with a question with choices, the corresponding prediction of the KIT model, and two rationales: LLMgenerated (from our pipeline) and crowdworkerwritten. We then ask them to make a preferential selection among the two rationales (see interface details in Appendix E.1.) We find low-tomoderate annotator agreement – Krippendorff's $\alpha = 0.13$ (Krippendorff, 2011) — for this study, indicating the subjective nature of the task. Related work (Wiegreffe et al., 2022) reported similar agreement statistics ($\alpha \in [0.05, 0.20]$) on comparison between LLM-generated and ECQA rationales.

Fine-grained comparison. Besides head-to-head comparison, we ask several 7-point Likert scale questions — adapted from prior work (Aggarwal et al., 2021; Wiegreffe et al., 2022) — targeted at comparing fine-grained aspects of both rationales. These aspects include: *sufficiency* in justifying the model's choice; *conciseness* (*i.e.*, degree of redundancy); *understandability*; *factuality* (*i.e.*, factual correctnes); *supportiveness* (*i.e.*, the degree to which the model prediction is supported); *refutation convincingness* (*i.e.*, the degree to which the



Figure 3: Distribution of fine-grained metrics between crowdworker (ECQA) and LLM-generated rationales — LLM-generated rationales were preferred over ECQA on majority of the metrics except conciseness.

unselected choices are convincingly refuted); *in-sightfulness* (*i.e.*, how much new information is captured.) New information can be new facts or reasoning not stated in the question and answer choices and potentially grounded on the knowl-edge evidence. We report the agreement statistics on individual aspects in the Appendix B.2.

310

311

312

313

314

315

316

317

318

319

321

322

323

324

325

326

327

330

331

332

333

334

335

336

337

341

343

344

4.2 Higher Preference of LLM Generations

Surprisingly, LLM-generated rationales were more frequently preferred (67.2% times) over crowdworker-written rationales (37.8% times.)The result showcases an improvement over previous work on generating corroboration only (no refutation) rationales (Wiegreffe et al., 2022) — 45.7%preference to LLM generations. The crowdworkerwritten ECQA rationales potentially outperformed those corroboration rationales on dimensions such as refutation convincingness, sufficiency, and sup-Moreover, our pipeline enabled portiveness. knowledge-guided rationale generation, whereas prior LLM-generated rationales (Wiegreffe et al., 2022) lacked such grounding and were abstractive. However, the LLMs in both studies differed, with our study employing a newer version (GPT-3.5) than the GPT-3 model used in their work. While some of the performance gain can be attributed to such model upgrades, we demonstrate via finegrained analysis how aspects of our rationale construct are correlated with crowd-worker preference.

4.3 Fine-grained Comparison

As shown in Figure 3, overall, crowd workers exhibited more preference for LLM-generated rationales over crowdworker-written ones on aspects such as insightfulness (*i.e.*, new information), refutation convincingness, and sufficiency. In fact, up to 80.4% of the LLM-generated rationales presented to the crowdworkers contained at least one statement grounded on external knowledge, thereby contributing to insightfulness (see Appendix D.1 for details.) Moreover, the refutation argument anchored on the topic of the question enabled a more convincing refutation. Therefore, the resulting LLM-generated rationales were sufficient to justify the model's choice for the QA task.

345

346

351

354

359

362

363

370

371

373

374

375

376

386

The preference for LLM-generated and crowdworker-written rationales was comparable for other aspects such as factuality, supportiveness, and understandability. However, crowd-workers rated LLM-generated rationales as more redundant which is unsurprising, given the tendency of the LLMs to generate verbose text.

Metrics	LLM-generated Preferred	Crowdworker-written Preferred
Factuality	0.29	0.04
Insightfulness	0.21	0.12
Conciseness	0.08	0.02
Convincingness	0.29	0.17
Sufficiency	0.28	0.14
Supportiveness	0.27	0.03
Understandability	0.27	0.01

Table 1: Spearman correlation between crowdworker preference of rationales — *weak* correlations are observed with p < 0.01 (strong statistical significance.)

Correlation to rationale preference. To understand what factors are important for the preference judgment, we compute Spearman correlation (Spearman, 1987) between the binary preference of both rationale types — *i.e.*, LLMgenerated and crowdworker-written - and the finegrained aspects (see Table 1.) The conciseness aspect lacked any correlation with either rationale type. Surprisingly, crowd-workers' preference for crowdworker-written rationales lacked any correlation with several other aspects, such as factuality, supportiveness, and understandability, while showcasing a very weak correlation with the rest. However, these fine-grained aspects exhibited a comparatively stronger positive correlation with the LLMgenerated rationales. Further analysis showcases that even when crowd-workers preferred ECQA rationales in a head-to-head comparison, LLMgenerated and crowdworker-written rationales exhibited almost similar ratings in the majority of the fine-grained aspects (see Appendix D.2.) Overall, the results indicate that human preference for LLMgenerated rationale can be captured by factoring in different fine-grained aspects, which can inform the design of automated mechanisms for estimating the suitability of a rationale for end-users.

5 Acceptability of LLM Rationalization

387

388

389

390

391

392

393

394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

While pairwise evaluations of preferences provided perspective on the relative quality of the rationales, we conducted another study to independently measure the acceptability of the LLM-generated rationales and collect absolute crowd-worker judgments across several aspects related to rationale quality. We evaluated rationales for both CSQA and OBQA dataset task to understand how generalizable these observations are. The key takeaways from the study are as follows:

- the overall acceptability of the rationales remained high similar to the comparative study.
- however, task and domain variation impacted the quality of the generated rationales.

5.1 Study Setting

In each of the 250 HITs per dataset (three different judges per HIT), a crowd-worker was presented with a question with choices, the corresponding prediction of the KIT model, and an LLM-generated rationale. Besides asking 7-point Likert scale questions on fine-grained aspects of a rationale - similar to the first study in Section 4 — we include two additional surface-level aspects: readability, i.e., the clarity of the provided justifications and gram*maticality*, adherence to grammatical rules. Finally, we ask for an overall judgment on quality, *i.e.*, the overall acceptability of a rationale (see interface details in Appendix E.1.) We again find low-tomoderate agreement – Krippendorff's $\alpha = 0.12$ for CSQA and 0.15 for OBQA dataset. Related work (Wiegreffe et al., 2022) reported slightly better agreement statistics ($\alpha = 0.28$) on the CSQA dataset (see Appendix B.2 for details.)

5.2 Favorability Towards LLM generations

On the overall acceptability metric, the LLMgenerated rationales received a notably positive rating from the participants for both CSQA ($\mu =$ $5.83, \sigma = 1.27$) and OBQA ($\mu = 5.89, \sigma = 1.50$). These independent observations reaffirm earlier takeaways (§ 4) and underscore that the LLMgenerated rationales of KIT models were viewed favorably by crowd-workers.

Fine-grained observations. As shown in Figure 4, for the newly introduced surface-level metric, readability, and grammaticality, the LLM-generated rationales received higher ratings in keeping with the previous work. In fact, for both datasets, for all of the richer aspects except *insightfulness* and



Figure 4: Crowdworkers' ratings showed similar distrbution for all metrics except insightfulness and concisenes. These metrics were rated lower for the more subjective CSQA dataset compared to the objective and scientific OBQA dataset.

conciseness, the ratings received were similar, *i.e.*, more positively rated. While the insightfulness metric was rated positively for OBQA, the rating was neutral to slightly negative for CSQA. Surprisingly, conciseness (*i.e.*, less redundancy) was rated positively for OBQA, whereas CSQA rationales were deemed more redundant, similar to the previous study. A plausible explanation for this discrepancy is the inherent subjectivity in CSQA (Wiegreffe et al., 2022), which can result in varying expectations regarding the information provided in the rationales. In contrast, the OBQA dataset is grounded in objective scientific facts, eliminating such subjectivity and leading to more consistent expectations among crowd-workers.

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

Metrics	Correlation CSQA	Correlation OBQA
Factuality	0.65	0.73
Insightfulness	0.38	0.67
Conciseness	0.09	0.6
Convincingness	0.70	0.80
Sufficiency	0.76	0.80
Supportiveness	0.54	0.76
Understandability	0.63	0.71
Readability	0.5	0.74
Grammar	0.33	0.62

Table 2: Spearman correlation between acceptability and the fine-grained aspects of a rationale — moderate to fairly strong correlation were observed with **strong statistical significance** (p < 0.01).

Correlation to overall acceptability. To understand what factors are important for the overall *acceptability* judgement, we compute Spearman correlation (Spearman, 1987) between *acceptability* and the fine-grained aspects (see Table 2.) For both the datasets, all aspects except *conciseness* show similar patterns — moderate to fairly strong positive correlation with acceptability. However, the rationales for the CSQA dataset (more subjective) exhibited a weaker correlation than the OBQA dataset rationales (more objective) in several aspects, such as conciseness, insightfulness, readability, and grammaticality. **Overall**, *the results indicate that human preference for rationale is more nuanced and can only be holistically captured by considering different fine-grained aspects. However, the quality of the generated rationale may vary depending on the task and domain and, consequently, impact human-preference judgment. Therefore, there is room for improvement in making generated rationales invariant to task and domain.*

460

461

462

463

464

465

466

467

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

495

496

497

498

499

500

501

502

503

504

505

506

6 Towards Credible Rationalization

In the earlier studies, similar to existing work (Aggarwal et al., 2021; Wiegreffe et al., 2022; Marasovic et al., 2022; Kayser et al., 2021), we evaluate LLM-generated rationales for cases where model prediction matches the ground truth. We now investigate the implications of rationalization without accounting for model errors, *i.e.*, faithful rationalization, and potential intervention strategies. Following are the key highlights of the study:

- rationalizing incorrect predictions drastically reduces human's trust in the LLM rationalizer.
- even lightweight guardrails can help intervene more than half of the incorrect predictions.

6.1 Trustworthiness of Generated Rationales

The reported accuracy of KIT models widely vary -64%-89.4% for CSQA² and 60.4%-89.6% for OBQA³. The reported human accuracy for the CSQA and OBQA datasets are 88.9% and 91.7%, respectively. Even as humans rationalize, the credibility of the rationalizer may diminish if they attempt to justify any incorrect decisions. Existing work on trust in explainable AI (XAI) literature (Hoff and Bashir, 2015; Schaefer et al., 2016; Stites et al., 2021; Smith-Renner et al., 2020) demonstrates that end-users' trust in a system degrades when encountering errors they can easily recognize due to familiarity and prior experience in a domain. Since the knowledge source for the CSQA and OBQA datasets is ConceptNet (Speer et al., 2017), a commonsense knowledge graph, humans are expected to have higher confidence about their knowledge in the domain. However, existing explainable NLP literature lack studies that investigate the relationship between model accuracy and

²https://www.tau-nlp.sites.tau.ac.il/

³https://leaderboard.allenai.org/open_book_qa/

humans' degree of trust in the context of free-text
rationales. Therefore, we replicate the study design of exploring trust in explanations for classification models (Stites et al., 2021) to confirm whether
the observations hold for knowledge-intensive QA
tasks in the commonsense domain.

Study design. We conducted a between-subject 513 study involving 22 participants (15 male and 7 fe-514 male) exploring two conditions: 66% (11 partici-515 pants) and 90% (11 participants) model accuracy. 516 The accuracy conditions reflect the two extremities 517 of existing knowledge-intensive task models (Yasunaga et al., 2021; Feng et al., 2020; Yasunaga 519 et al., 2022). The study consisted of three phases: an introduction to the study, a quiz phase, and a 521 follow-up survey. In the quiz phase, the participants answered 15 QA tasks. The 15 tasks were randomly selected from the CSQA (8 QAs) and 524 OBQA (7 QAs) datasets. Depending on the study conditions, for X% of those N questions, where $X \in \{66, 90\}$, the KIT model made accurate pre-527 dictions, and the rest of the predictions were inaccurate. The KIT model prediction and LLMgenerated rationale of a QA task were revealed after a participant submitted their response to avoid 531 bias. Then, the participants were asked whether they agreed with the model prediction and had to 533 rate their impression of the rationale on a scale of 1 to 7 (1 = actively misleading and 7 = helpful.) 535 After the quiz phase, the participants completed a survey adapted from the Trust Scale recommended for XAI (Hoffman et al., 2018). The survey con-538 tained questions that asked participants to rate several aspects related to the quiz phase tasks, such as 540 the agreement with rationales and the participants' trust and reliance on the LLM-generated rationale. 542 All of these required participants to work slowly enough to be able to read all the items, thereby making the studies long-running and rather un-545 suitable for crowd platforms according to existing work (Douglas et al., 2023). Therefore, we opted 547 for internal recruitment as an additional quality control mechanism, inviting participants internally via a Slack campaign at Company X. None of the par-550 ticipants are authors of the paper (see Appendix C.)

6.2 Confirmatory Study Results

552

556

The agreement statistics of the participants reflect both the study conditions — 67.27% and 86.07% for lower and higher accuracy models, respectively. Figure 5a summarizes the participants' impression



(a) Impact on user perception (b) XAI Trust Scale feedback

Figure 5: (a) Irrespective of agreement or disagreement with the KIT model prediction, participants indicated a more negative impression about the rationalization of the lower confidence model prediction. (b) Participant feedback on trust scale indicates lower confidence for lower accuracy model rationalization.

557

558

559

560

561

562

563

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

584

585

586

588

589

590

592

of a rationale immediately after viewing the model prediction. When the participants disagreed with the model prediction, they exhibited a stronger negative impression about the rationales for the 66%accuracy condition compared to the 90% accuracy condition. Even when participants agreed with the model prediction, their impression of the rationales remained more negative. Our intuition is that the higher disagreement with the model coupled with observing the faithful rationalization of the incorrect prediction negatively impacted participants' perception of the reliability of the rationales. We confirm these observations by analyzing the results of the follow-up survey (see Figure 5b.) Unsurprisingly, participants for the 66% accuracy condition rated their confidence in the generated rationales and the reliability of the rationalizer significantly lower compared to the 90% accuracy condition. The trends in Figure 5 are observed with strong statistical significance, except for participant feedback on satisfaction with rationale (see Appendix C.2.)

6.3 A Review-then-Rationalize Framework

Motivated by the observations from the preliminary study, we create a two-stage review-thenrationalize (see Figure 6) pipeline to evaluate the impact of intervening incorrect model predictions before rationalization. The pipeline instruments a *reviewer* module that employs another model (GPT-3.5 text-davinci-003 (temperature = 0)) to evaluate the correctness of the knowledge-intensive task model and refrain from rationalizing potentially incorrect decisions.

We opted for LLMs as reviewers due to their reported proficiency in natural language understanding. Depending on the task and data domain, the suitability of the reviewer model may vary. Given



Figure 6: Self-consistency-based Reviewer-intervene for any disagreement with the KIT model prediction.

594

599

607

610

613

617

624

the complexity of knowledge-intensive tasks, we employ a self-consistency-based decoding strategy (Wang et al., 2022) as opposed to greedydecoding to ensure robustness. More specifically, we independently pose the same QA task N (=5)times to the reviewer and select the final response via majority voting. The reviewer then compares the model's prediction with its prediction and activates the rationalizer only when both models agree. A cookie-cutter rationale or no rationale may be communicated to the end-user in a disagreement.

Detect	Prediction Errors	Errors Intervened		
Dataset	(Test Set)	Greedy Decoding	Self-consistency	
CSQA	321	166 (51.71%)	187 (58.26 %)	
OBQA	155	102 (65.81%)	$110({\bf 70.97\%})$	

Table 3: The review-then-rationalize pipeline helps intervene incorrect predictions of a knowledge-intensive task model. The self-consistency-based reviewer outperforms the greedy decoding-based reviewer.

As shown in Table 3, for knowledge-intensive tasks such as CSQA and OBQA, the proposed pipeline helps intervene up to 58% and 71% of the incorrect predictions. Unsurprisingly, the selfconsistency-based reviewer outperforms the greedy decoding-based reviewer. Overall, the results draw attention to the importance of responsibly communicating LLM-generated rationales to humans and consequently, instrumenting guardrails as an effective intervention strategy.

Related Work 7

Free-text Rationale Generation. Existing works 615 highlight the effectiveness of free-text rationales 616 in justifying a model's decision to humans in vision (Hendricks et al., 2016; Park et al., 2018) and 618 text domains (Camburu et al., 2018; Ehsan et al., 2018; Narang et al., 2020). Due to cost and generalizability implications of supervised rationale generation, we employ few-shot prompting to elicit rationales from LLMs following existing work (Wiegreffe et al., 2022; Marasović et al., 2021). Both these approaches generate abstractive, corroborative, and faithful rationales. In contrast, we explore the generation of knowledge-guided, corroborative and refutation complete, and credible rationales.

Guided text generation. Developing approaches to avoid hallucinations and factual inaccuracies in LLM-generated text is a new area of research. Retrieval augmented generation (RAG) infuses external knowledge (Peng et al., 2023; Lazaridou et al., 2022), such as knowledge-bases and web documents, while prompting LLMs to help generate responses. We employ a similar strategy during rationalization by conditioning the LLM generation on the retrieved evidence for a given task.

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

Credible text generation. Studies in explainable AI literature (Smith-Renner et al., 2020; Hoff and Bashir, 2015; Schaefer et al., 2016; Stites et al., 2021) demonstrate that for low-quality models, providing faithful explanations further degraded user's trust. Unlike existing work on free-text explanation (Wiegreffe et al., 2022; Marasović et al., 2021), we explore how end-users' trust may be impacted by faithful rationalization of varying degrees of incorrect model predictions. ReXC (Majumder et al., 2021) augments rationales - generated in a self-rationalization framework --- with background knowledge to improve a model's task performance, such as natural language inference and visual commonsense reasoning. To rectify incorrect LLM responses, identified via a self-consistency-based intervention approach, the Verify-then-Edit framework (Zhao et al., 2023) leverages external knowledge to repair reasoning chains of the corresponding chain-of-thought prompts. FARM (Mei et al., 2023) utilizes trustworthy external sources within a predict-then-generate framework that aims to intervene in harmful content generation using LLMs. To credibly rationalize KIT model predictions, we explore a review-then-rationalize framework where a self-consistency-based reviewing approach identifies potential prediction inaccuracies and ensures credible rationale generation.

8 Conclusion

We evaluate LLMs' capacity to generate effective rationales for knowledge-intensive tasks in a fewshot knowledge-guided setting. We additionally investigate the implications of employing LLMs as rationalizers of an imperfect model and highlight the negative impact on users' trust. Observations from our studies highlight room for improvement in aspects such as task and domain invariant rationalization and robust intervention strategies for real-world usage.

697

698

701

702

704

708

710

711

713

714

716

718

719

721

723

725

726

9 Limitations

Scrutinizing LLM-generated rationales. While external knowledge-guided generation offers promise (Peng et al., 2023; Mallen et al., 2023), LLM-generated rationales may still suffer from hallucinations. Our experiments highlight that the LLM-generated rationale is not entirely grounded on retrieved knowledge. Even though crowdworkers positively rated the factuality and insightfulness of the generated rationales, additional scrutiny is required before deploying such rationalizers in mission-critical tasks. To this end, the review-then-rationalize framework may be expanded to further scrutinize the rationales by adopting recent work on an LLM's factual knowledge measurement (Pezeshkpour, 2023; Dong et al., 2023) and hallucination identification (Manakul et al., 2023; Elaraby et al., 2023; Mündler et al., 2023) and reduction (Zhao et al., 2023; Mei et al., 2023), and explainable evaluation (Xu et al., 2023). Fairwashing vs. credible rationalization. The accuracy of our self-consistency-based reviewer can be further improved to intervene in a higher proportion of incorrect KIT model predictions. However, critiques of XAI tools have raised concerns about fairwashing, i.e., misleading users into trusting biased or incorrect models (Alikhademi et al., 2021). For example, simply averting potential faithful yet incorrect rationalization, identified by the reviewer, may increase end-users' trust due to an illusion of a highly performant rationalizer (Aïvodji et al., 2019). Such fairwashing may have catastrophic consequences if employed in real-world applications such as in the medical domain, hiring platforms, and credit agencies. Recent work (Alikhademi et al., 2021) proposes a framework for evaluating XAI tools with respect to their capabilities for detecting and addressing issues of bias and fairness as well as their capacity to communicate these results to their users clearly. Therefore, future implementations of the credible rationale should adopt similar strategies to safeguard against such issue. Future work may explore different communication strategies during prediction errors, such as communicating the disagreement to the experts-in-the-loop, providing rationales with a disclaimer, and employing stronger reviewers to repair the prediction on the fly and then rationalize, among others.

727 Scaling responsibly. An often overlooked aspect728 of the recent popularity of LLMs has been *Green*

AI (Schwartz et al., 2020). When the ML deployment pipeline is considered as a whole, inference consumes most compute resources, accounting for anything between 70% to 90% (Weng et al., 2022; Wu et al., 2022). Knowledge distillation approaches can be adopted to avoid costly pretraining (Wang et al., 2023). Furthermore, materialization of rationales to avoid repeating rationalizing the same task can be possible approaches to handle such issues.

729

730

731

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

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

Crowdsourcing study constraints. As we conducted the crowdsourced study on Amazon Mechanical Turk, our findings may not generalize to other platforms and feedback provided in in-person lab-based studies. Moreover, we observed low agreement among the annotators — similar to prior work (Wiegreffe et al., 2022) - due to the subjectivity of the QA tasks. Future work may explore conducting large-scale studies with better quality control mechanisms (such as hiring private firms with dedicated teams similar to (Aggarwal et al., 2021) and conducting in-house studies with experts. Such a setting also makes to possible to collect additional insights into the thought process of the participants. However, conducting such large-scale studies in an in-person setup introduces time and logistics constraints. To this end, recent LLM-based reference-free approaches (Liu et al., 2023) to scale-up evaluation offers promise. Although, it is unclear whether such evaluation strategies apply to subjective metrics of rationale quality studied in our work. Therefore, future studies may explore how such reference-free judgements align with human judgements similar to (Pezeshkpour, 2023).

References

- Shourya Aggarwal, Divyanshu Mandowara, Vishwajeet Agrawal, Dinesh Khandelwal, Parag Singla, and Dinesh Garg. 2021. Explanations for CommonsenseQA: New Dataset and Models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Online. Association for Computational Linguistics.
- Ulrich Aïvodji, Hiromi Arai, Olivier Fortineau, Sébastien Gambs, Satoshi Hara, and Alain Tapp. 2019. Fairwashing: the risk of rationalization. In *International Conference on Machine Learning*, pages 161–170. PMLR.
- Kiana Alikhademi, Brianna Richardson, Emma Drobina, and Juan E Gilbert. 2021. Can explainable ai explain

836

837

unfairness? a framework for evaluating explainable ai. *arXiv preprint arXiv:2106.07483*.

780

781

782

790

794

795

796

797

799

800

805

810

811

812

815

817

819

821

822

823

824

825

826

827

830

832

835

- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- 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.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. *Advances in Neural Information Processing Systems*, 31.
- Miruna-Adriana Clinciu, Arash Eshghi, and Helen Hastie. 2021. A study of automatic metrics for the evaluation of natural language explanations. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2376–2387, Online. Association for Computational Linguistics.
- Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Zhifang Sui, and Lei Li. 2023. Statistical knowledge assessment for large language models. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Benjamin D Douglas, Patrick J Ewell, and Markus Brauer. 2023. Data quality in online humansubjects research: Comparisons between mturk, prolific, cloudresearch, qualtrics, and sona. *Plos one*, 18(3):e0279720.
- Upol Ehsan, Brent Harrison, Larry Chan, and Mark O Riedl. 2018. Rationalization: A neural machine translation approach to generating natural language explanations. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 81–87.
- Mohamed Elaraby, Mengyin Lu, Jacob Dunn, Xueying Zhang, Yu Wang, and Shizhu Liu. 2023. Halo: Estimation and reduction of hallucinations in opensource weak large language models. *arXiv preprint arXiv:2308.11764.*
- Yanlin Feng, Xinyue Chen, Bill Yuchen Lin, Peifeng Wang, Jun Yan, and Xiang Ren. 2020. Scalable multihop relational reasoning for knowledge-aware question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1295–1309, Online. Association for Computational Linguistics.
- Sai Gurrapu, Ajay Kulkarni, Lifu Huang, Ismini Lourentzou, Laura Freeman, and Feras A Batarseh.

2023. Rationalization for explainable nlp: A survey. *arXiv preprint arXiv:2301.08912*.

- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.
- Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. 2016. Generating visual explanations. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 3–19. Springer.
- Kevin Anthony Hoff and Masooda Bashir. 2015. Trust in automation: Integrating empirical evidence on factors that influence trust. *Human factors*, 57(3):407– 434.
- Robert R Hoffman, Shane T Mueller, Gary Klein, and Jordan Litman. 2018. Metrics for explainable ai: Challenges and prospects. *arXiv preprint arXiv:1812.04608*.
- Shiyuan Huang, Siddarth Mamidanna, Shreedhar Jangam, Yilun Zhou, and Leilani H Gilpin. 2023. Can large language models explain themselves? a study of llm-generated self-explanations. *arXiv* preprint arXiv:2310.11207.
- Maxime Kayser, Oana-Maria Camburu, Leonard Salewski, Cornelius Emde, Virginie Do, Zeynep Akata, and Thomas Lukasiewicz. 2021. e-vil: A dataset and benchmark for natural language explanations in vision-language tasks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1244–1254.
- John King and Roger Magoulas. 2015. 2015 data science salary survey. O'Reilly Media, Incorporated, Sebastopol, CA, USA.
- Klaus Krippendorff. 2011. Computing krippendorff's alpha-reliability.
- Angeliki Lazaridou, Elena Gribovskaya, Wojciech Stokowiec, and Nikolai Grigorev. 2022. Internetaugmented language models through few-shot prompting for open-domain question answering. *arXiv preprint arXiv:2203.05115*.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.

Andreas Madsen, Siva Reddy, and Sarath Chandar. 2022. Post-hoc interpretability for neural nlp: A survey. *ACM Computing Surveys*, 55(8):1–42.

893

895

897

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

924

925

932

934

937

938

939

941

942

- Bodhisattwa Prasad Majumder, Oana-Maria Camburu, Thomas Lukasiewicz, and Julian McAuley. 2021.
 Knowledge-grounded self-rationalization via extractive and natural language explanations. arXiv preprint arXiv:2106.13876.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023.
 When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.
 - Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896*.
- Ana Marasovic, Iz Beltagy, Doug Downey, and Matthew Peters. 2022. Few-shot self-rationalization with natural language prompts. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 410–424, Seattle, United States. Association for Computational Linguistics.
 - Ana Marasović, Iz Beltagy, Doug Downey, and Matthew E Peters. 2021. Few-shot selfrationalization with natural language prompts. *arXiv preprint arXiv:2111.08284*.
 - Alex Mei, Sharon Levy, and William Yang Wang. 2023. Foveate, attribute, and rationalize: Towards physically safe and trustworthy ai. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11021–11036.
 - Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*.
- Niels Mündler, Jingxuan He, Slobodan Jenko, and Martin Vechev. 2023. Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation. *arXiv preprint arXiv:2305.15852*.
- Sharan Narang, Colin Raffel, Katherine Lee, Adam Roberts, Noah Fiedel, and Karishma Malkan. 2020. Wt5?! training text-to-text models to explain their predictions. arXiv preprint arXiv:2004.14546.
- Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. 2018. Multimodal explanations: Justifying decisions and pointing to the evidence. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 8779–8788.

Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, et al. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813*. 943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

- Pouya Pezeshkpour. 2023. Measuring and modifying factual knowledge in large language models. *arXiv* preprint arXiv:2306.06264.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4932–4942, Florence, Italy. Association for Computational Linguistics.
- Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585.
- Emily Reif, Daphne Ippolito, Ann Yuan, Andy Coenen, Chris Callison-Burch, and Jason Wei. 2021. A recipe for arbitrary text style transfer with large language models. *arXiv preprint arXiv:2109.03910*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Kristin E Schaefer, Jessie YC Chen, James L Szalma, and Peter A Hancock. 2016. A meta-analysis of factors influencing the development of trust in automation: Implications for understanding autonomy in future systems. *Human factors*, 58(3):377–400.
- Timo Schick and Hinrich Schütze. 2020. Exploiting cloze questions for few shot text classification and natural language inference. *arXiv preprint arXiv:2001.07676.*
- Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. 2020. Green ai. *Communications of the ACM*, 63(12):54–63.
- Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. 2020. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*.

999

- 1014
- 1015
- 1017 1018
- 1019
- 1021 1022
- 1023 1024
- 1025
- 1026 1027
- 1028 1029
- 1030 1031
- 1032
- 1034

1035

1038

- 1041

1045 1046 1047

1049

1050 1051

1048

Alison Smith-Renner, Ron Fan, Melissa Birchfield, Tongshuang Wu, Jordan Boyd-Graber, Daniel S Weld, and Leah Findlater. 2020. No explainability without accountability: An empirical study of explanations and feedback in interactive ml. In Proceedings of the 2020 chi conference on human factors in computing systems, pages 1–13.

- Charles Spearman. 1987. The proof and measurement of association between two things. The American journal of psychology, 100(3/4):441-471.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In Proceedings of the AAAI conference on artificial intelligence, volume 31.
- Mallory C Stites, Megan Nyre-Yu, Blake Moss, Charles Smutz, and Michael R Smith. 2021. Sage advice? the impacts of explanations for machine learning models on human decision-making in spam detection. In International Conference on Human-Computer Interaction, pages 269-284. Springer.
- Jiao Sun, Swabha Swayamdipta, Jonathan May, and Xuezhe Ma. 2022. Investigating the benefits of freeform rationales. In Findings of the Association for Computational Linguistics: EMNLP 2022, pages 5867–5882, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseOA: A guestion answering challenge targeting commonsense knowledge. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4149-4158, Minneapolis, Minnesota. Association for Computational Linguistics.

- Chenhao Tan. 2021. On the diversity and limits of human explanations. arXiv preprint arXiv:2106.11988.
- Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. 2023. Scott: Self-consistent chain-of-thought distillation. arXiv preprint arXiv:2305.01879.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171.
- Samuel F Way, Daniel B Larremore, and Aaron Clauset. 2016. Gender, productivity, and prestige in computer science faculty hiring networks. In Proceedings of the 25th International Conference on World Wide Web, pages 1169–1179, New York, NY, USA. ACM.
- Qizhen Weng, Wencong Xiao, Yinghao Yu, Wei Wang, Cheng Wang, Jian He, Yong Li, Liping Zhang, Wei

Lin, and Yu Ding. 2022. {MLaaS} in the wild: Workload analysis and scheduling in {Large-Scale} heterogeneous {GPU} clusters. In 19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22), pages 945–960.

1052

1053

1055

1057

1058

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

- Sarah Wiegreffe, Jack Hessel, Swabha Swayamdipta, Mark Riedl, and Yejin Choi. 2022. Reframing human-AI collaboration for generating free-text explanations. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 632-658, Seattle, United States. Association for Computational Linguistics.
- Sarah Wiegreffe and Ana Marasović. 2021. Teach me to explain: A review of datasets for explainable natural language processing. arXiv preprint arXiv:2102.12060.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.
- Carole-Jean Wu, Ramya Raghavendra, Udit Gupta, Bilge Acun, Newsha Ardalani, Kiwan Maeng, Gloria Chang, Fiona Aga, Jinshi Huang, Charles Bai, et al. 2022. Sustainable ai: Environmental implications, challenges and opportunities. Proceedings of Machine Learning and Systems, 4:795–813.
- Yunshu Wu, Hayate Iso, Pouya Pezeshkpour, Nikita Bhutani, and Estevam Hruschka. 2023. Less is more for long document summary evaluation by llms. arXiv preprint arXiv:2309.07382.
- Wenda Xu, Danqing Wang, Liangming Pan, Zhenqiao Song, Markus Freitag, William Yang Wang, and Lei Li. 2023. Instructscore: Towards explainable text generation evaluation with automatic feedback. arXiv preprint arXiv:2305.14282.
- Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren, Xikun Zhang, Christopher D Manning, Percy S Liang, and Jure Leskovec. 2022. Deep bidirectional language-knowledge graph pretraining. Advances in Neural Information Processing Systems, 35:37309– 37323.
- Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, 1099 Percy Liang, and Jure Leskovec. 2021. QA-GNN: 1100 Reasoning with language models and knowledge 1101 graphs for question answering. In Proceedings of 1102 the 2021 Conference of the North American Chapter 1103 of the Association for Computational Linguistics: Hu-1104 man Language Technologies, pages 535-546, Online. 1105 Association for Computational Linguistics. 1106

- 1107 1108 1109
- 1110
- 1111 1112

1114

1115

1116

1117

1118

1119

1138

1139

1140

- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei Qin, and Lidong Bing. 2023. Verify-and-edit: A knowledge-enhanced chain-of-thought framework. *arXiv preprint arXiv:2305.03268*.

A Prompts and Rationales

In this section, we provide additional details regarding the prompts corresponding to the faithful and credible rationalization workflows.

A.1 Faithful Rationalization

Figure 7 outlines the few-shot prompt structure. 1120 Table 4 elaborates on the prompt design shown 1121 in (Figure 2 and Figure 7). Each example in the 1122 few shot prompt includes the question and answer 1123 choices, the KIT model selected answer, the knowl-1124 edge facts extracted from ConceptNet for each 1125 choice, and the expert-written question topic and 1126 rationale that act as input to GPT-3.5 text davinci 1127 003. While we show only two few-shot examples, 1128 in practice, we use five examples per prompt. As 1129 explained in the Section 2, due the token limit im-1130 posed by the GPT-3.5 API, we can include from 1131 5-8 examples depending on the length of the knowl-1132 edge facts. Given the prompt, *i.e.*, examples fol-1133 lowed by an unseen question and answer choices, 1134 KIT model selected answer, and extracted knowl-1135 edge, the LLM greedily generates the question 1136 topic and the rationale for the model prediction. 1137



Figure 7: An example in the few-shot prompt: the QA and External Knowledge components are retrieved and the topic and the rationale are expert authored.

To design the initial prompt, we take inspiration from existing work (Wiegreffe et al., 2022; Peng et al., 2023; Lazaridou et al., 2022; Zhao et al., 2023) to experiment with the prompt layout. We 1141 experimented with approximately 6 different lay-1142 outs in the OpenAI playground ⁴ using 10 train-1143 ing examples on the CSQA and OBQA datasets. 1144 In deciding the number of few-shot examples, we 1145 considered the maximum context window size of 1146 GPT-3.5 text-davinci-003, which is 4097 tokens. 1147 We observed that depending on the datasets and the 1148 length of the factual statements retrieved from Con-1149 ceptNet, five to eight few-shot examples fit into the 1150 token constraints. After finalizing the prompt lay-1151 out, we developed a pool of 40 expert-written (i.e., 1152 authors of these papers) examples. We randomly 1153 selected 5 expert-written examples for each test 1154 instance to ensure uniformity across datasets and 1155 instances. Similar to prior work (Wiegreffe et al., 1156 2022), we focused on developing a general few-1157 shot prompting strategy for generating knowledge-1158 enhanced and refutation complete rationale that 1159 could be successful when no additional (large) val-1160 idation set for parameter tuning is available. We 1161 prompt the LLM to generate a topic of the question 1162 and a rationale similar to the provided few-shot 1163 examples. Therefore, our approach explicitly con-1164 ditions the rationale generation on the question 1165 topic and the knowledge facts. FARM (Mei et al., 1166 2023) employs a similar topic-focused generation 1167 for question answering. Given a question, the LLM 1168 is initially prompted to generate a question con-1169 text — augmented with information retrieved from 1170 trustworthy sources — to generate a safe response. 1171 Such strategies have been shown to be very effec-1172 tive (Radford et al., 2019; Brown et al., 2020; Shin 1173 et al., 2020; Schick and Schütze, 2020), even in 1174 complex generation tasks (Reif et al., 2021). 1175

Relevance to ECQA rationales. The pipeline for ECQA (Aggarwal et al., 2021) rationale generation and the knowledge-guided LLM rationalization have several similarities. As shown in Table 4, the rationalization pipeline provides conceptnet assertions corresponding to the selected answer and rejected choices as context within the prompt. ECQA crowdsourcing pipeline also prompted crowdworkers to use the positive facts about the selected answer and negative facts about the other choices as guides to craft the eventual free-flow explanation. Therefore, in both cases, the rationalizer pipeline, be it crowdworker or LLMs, were knowledgeguided. However, in case of LLMs the source of knowledge guidance is external, *i.e.*, Concept-

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

⁴https://platform.openai.com/playground

Question: At the end of your meal what will a waiter do Choices: A. serve food B. eat C. set table D. serve meal E. present bill Selected answer: E. present bill Knowledge for present bill: [waiter can typically do present bill, bill is generally created by waiter, restaurant generally causes bill, ...] Knowledge for set table: [waiter can typically do set table, ...] Knowledge for serve food: [waiter can typically do serve food, Knowledge for serve meal: [waiter can typically do serve meal, . . .] Knowledge for eat: [cook generally causes meal, . . .] The topic of the question and the corresponding explanation for the selected answer "present bill" are as follows: Topic: Restaurant Service after meal Why? Commonsense suggests that a waiter, who is generally located in a restaurant, typically presents a bill. Therefore, the answer is "present bill" because this is a common practice at the end of a meal in a restaurant Why not other options? While a waiter can serve food, set the table, and serve a meal, these actions typically occur before or during the meal, not at the end. The option 'eat' is not suitable as it is not a typical duty of a waiter during their service. _____ Question: He waited for his friend at the squash court, but he was worried his friend thought he meant the at the other end of the public what? Choices: A. country club B. rich person's house C. pool D. park E. fitness center Selected Answer: D. park Knowledge for park :[squash court is generally located in park, play is generally located in squash court, . . .] Knowledge for fitness center : [squash court is generally located in fitness center, ...] Knowledge for country club :[squash court is generally located in country club, . . .] Knowledge for pool :[...] Knowledge for rich person's house :[...] The topic of the question and the corresponding explanation for the selected answer "park" are as follows: topic: Public Spaces and Miscommunication Why? The answer is park because commonsense suggests that a squash court is generally located in a park. This implies that there could be another squash court at the other end of the park, leading to the friend's confusion. Why not other options? While a squash court can also be located in a fitness center or country club, these locations are not typically public spaces with multiple ends. A pool or a rich person's house are less likely to have multiple squash courts, making them less likely to be the source of the friend's confusion. ______ Question: What should the bean bag chair sit on? Choices: A. house B. den C. family room D. wood E. floor Selected Answer: E. floor Knowledge for present floor: [...] Knowledge for house: [...] Knowledge for den: [...] Knowledge for family room: [...] Knowledge for wood: [...]

The topic of the question and the corresponding explanation for the selected answer "present bill" are as follows:

Table 4: Example of a prompt with two training examples for CSQA and an unseen question for which the LLM generated a rationale. In practice, we provide five examples.

1191Net (Speer et al., 2017), whereas for ECQA, the1192crowdworkers themselves crafted the supporting1193facts before rationalizing.

1194 A.2 Credible Rationalization

Table 5 showcases the prompt design for the 1195 *Reviewer* model within the credible rationalizer 1196 pipeline (Figure 6. Each of the five examples in the few shot prompt includes the question and answer 1198 choices that act as input to GPT-3.5 text davinci 1199 003. In practice, we use five examples per prompt. 1200 Given the prompt, *i.e.*, examples followed by an unseen question and answer choices, the LLM greed-1203 ily generates a response, *i.e.*, predicts an answer from the choices. We repeat the process five times 1204 and select a response based on majority voting. We 1205 randomly sample five questions from the 40 expertwritten rationale pool as few-shot examples. 1207

Question: At the end of your meal what will a waiter do? Choices: A. serve food B. eat C. set table D. serve meal E. present bill Selected answer: E. present bill Question: He waited for his friend at the squash court, but he was worried his friend thought he meant the at the other end of the public what?

Choices: A. country club B. rich person's house C. pool D. park E. fitness center Selected Answer: D. park

Question: What should the bean bag chair sit on? Choices: A. house B. den C. family room D. wood E. floor Selected Answer:

Table 5: A prompt with two training examples for CSQA and an unseen question for the Reviewer to answer.

A.3 LLM-generated Rationales

Table 6 a few non-cherry picked examples of LLM-1209generated rationales. We show examples of ra-1210tionales for CSQA dataset questions generated by1211both LLM and humans, *i.e.*, crowdworkers in the1212ECQA dataset. Since there is no crowdsourced1213dataset of OBQA rationales, we only show LLM-1214generated rationales. Note that the LLM greedily1215

1216generates a topic of the question and a rationale1217with corroboration ("Why?") and refutation ("Why1218Not") components. However, as shown in Figure 2,1219the these two components are extracted from the1220generated output to construct the eventual ratio-1221nales (*i.e.*, the formatting step at the end.)

CSQA Dataset
Question: What should the bean bag chair sit on? Choices: A. house B. den C. family room D. wood E. floor Selected Answer: E. floor
LLM-generated rationale: The answer is floor because the common- sense knowledge clearly indicates that a bean bag chair is generally located in a floor. While a bean bag chair can be placed in a house, den, family room, or on wood, the floor is the most common place for a bean bag chair to be located.
ECQA rationale: Bean bag chair is a seat people sit on which is generally put on a floor. A bean bag chair should sit on a floor and not on anything else from the other options.
OBQA Dataset
Question: Rainbows are always found after what? Choices: A. A fire B. A tornado C. Rainfall D. Cereal Selected Answer: C. Rainfall
LLM-generated rationale: The answer is Rainfall because rainbows are always found after rain. This is because the sunlight is refracted by the raindrops in the air creating the rainbow. A fire, a tornado, and cereal do

raindrops in the air, creating the rainbow. A fire, a tornado, and cereal do not have any relation to rainbows.

Table 6: Rationales for CSQA generated by LLM and humans (ECQA) and for OBQA generated by LLM.

B Crowd Study Details

1222

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1239

1240

1241

1242

1243

1244

We provide more details regarding both crowdworker studies such as additional statistics related to the crowd study and quality control mechanisms.

B.1 Quality Control and Payment

In order to enforce quality throughout evaluation, we use a hidden built-in Javascript function to compute time per HIT spent and perform attention checks by inserting questions with specific instructions randomly within a HIT. We filter out any annotator who completed the tasks in an unreasonably low time, or failed their attention checks. To mitigate individual annotator bias, we also ensure that each experiment in a study has a substantial number of distinct crowdworkers. See Tables 7 and 8 for details regarding the inter-annotaror agreement for the comparison study. For both studies, we used a pay rate of USD 12.00/hr. We performed periodic check to ensure that the median HIT completion time remains commensurate to approximately the pay rate. Median times reported for the comparative study was 208 seconds (paid at 80 cents each) the acceptability study was 110 seconds (paid at

Approach	LLM-generated	ECQA
Factuality	0.07	0.05
Insightfulness	0.15	0.03
Conciseness	-0.04	-0.01
Convincingness	0.09	0.03
Sufficiency	0.08	0.07
Support	0.08	-0.01
Understandability	0.09	0.06
Preference	0.13	0.13

Table 7: Inter annotator agreement (Krippendorff's α) of crowdworkers on the fine-grained aspects of a rationale evaluated in the head-to-head comparison study.

40 cents each.) To ensure the quality of responses, 1245 we require annotators in Australia, New Zealand, 1246 United Kingdom, United States, and Canada as a 1247 proxy for English competency. We only selected 1248 workers with a past approval rate > 98% and who 1249 have completed over 5000 HITs. We documented a 1250 worker's HIT submission time and performed atten-1251 tion checks within each HIT to enforce quality con-1252 trol. Note that each crowd worker was presented 1253 with detailed instructions about the study interface 1254 and performed an example task as a warm-up. 1255

Dataset	CSQA	OBQA
Factual	0.02	0.03
Insightful	-0.06	-0.04
Concise	-0.15	-0.17
Convincing	0.08	0.13
Sufficient	0.07	0.08
Support	-0.012	-0.002
Understandable	0.02	0.04
Readability	-0.05	-0.02
Grammar	-0.15	-0.16
Acceptability	0.12	0.15

Table 8: Inter annotator agreement (Krippendorff's α) of crowdworkers on all the coarse- and fine-grained aspects of a rationale evaluated in the acceptability study.

B.2 Annotator Statistics

We now report the number of distinct crowd anno-1257 tators and the median and mean number of HITs 1258 completed for each experiment. For the head-to-1259 head comparison study, there were 750 HITs in to-1260 tal. There were 29 unique annotators with a median 1261 of 10 (mean = 21.86) HITs completed by an annotator. For the acceptability study, there 750 HITs 1263 for each of the two datasets CSQA and OBQA. For 1264 the CSQA dataset, there were 25 unique annotators 1265 with a median of 7 (mean = 28.80) HITs completed 1266 by an annotator. For the OBQA dataset, there were 1267 30 unique annotators with a median of 7 (mean 1268 = 25.00) HITs completed by an annotator. More 1269 detailed breakdowns of inter-annotator agreement 1270 for both studies are reported in Tables 7 and 8. 1271

1274

1275

1276

1277

1278

1279

1280

1281

1282

1284

1285

1286

1287

1288

1289

1290

1293

1294

1295

1296

1297

1298 1299

1300

1301

1303 1304

1305

1306

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

C Credible Rationalization Study

We now provide relevant information complementing the observations obtained in the preliminary study regarding credible rationalization.

C.1 Study Details

Participants. The participants of the preliminary study were all from Company X. However, we still performed attention checks in the preliminary study. The participants were unaware of the hypothesis and evaluation objective of the study. None of the participants are authors of the paper. Out of the 20 participants in the study, 15 were male and 5 were female. The representation of the female participants (25%) compares favorably with recent estimates of 15% women in tenure-track faculty in computing (Way et al., 2016) and 20%women in data science positions worldwide (King and Magoulas, 2015). One-fourth of the participants held a Bachelor degree and the rest completed graduate school or higher. Due to the complexity and longer duration of this study, we wanted to ensure the participation of higher quality participants by such selective recruitment.

Phases. We first collected participants' demographic information and then provided detailed instructions about the subsequent phases: a quiz phase consisting of a collection of tasks and a follow-up survey. The survey is adapted from the Trust Scale recommended for XAI (Hoffman et al., 2018). We opted for a follow-up survey rather than after each task completion following Hoffman et al. (Hoffman et al., 2018) — "the questions are appropriate for scaling after a period of use, rather than immediately after a rationale has been given." Besides questions related to the trust scale, we also asked participants to rate their overall acceptability of the rationales on a scale of 1 to 5. Note that the acceptability rating scale is different from the earlier studies in Section 5 and 4 to conform with the Trust Scale ratings (Hoffman et al., 2018).

C.2 Feedback Statistics

We conducted *Mann-Whitney U test* to measure the statistical significance of the differences between the 66% and 90% accurate model conditions, along various credibility metrics proposed in Section 6. The Mann-Whitney U test is a non-parametric test to measure the significance of difference in distribution of two independent sample, *i.e.*, accuracy conditions in this study. As shown in Table 9, participant feedback on individual task indicates a higher disagreement with lower confidence model prediction and a more negative impression regarding the corresponding rationale. The differences is significant both cases *i.e.*, when participants either agreed or disagree with the KIT model prediction. Table 10 reports the summary of participant feedback during the post-quiz survey — participants exhibited a more negative impression regarding the corresponding rationale. For all of the aspects except *statisfaction*, the difference in participant feedback between the accuracy conditions were statistically significant.

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1361

1362

1363

1364

1365

1366

1367

1369

Table 11 summarizes the observations from the quiz phase, *i.e.*, participant agreement statistics with the model prediction and participants' impression of the corresponding rationale. The agreement statistics (overall = 76.67%) of the participants reflect both the study conditions — 67.27% and 86.07%, respectively. Due to the subjective nature of the tasks, especially in the CSQA dataset, a few participants were unsure whether to agree or disagree with the model predictions, further reflecting the difficulty of the tasks.

D Additional Experiments and Analysis

We now present details of various user study observations, discussed briefly in earlier sections.

D.1 Degree of Knowledge Grounding

While our proposed knowledge-graph-based retrieval augmented LLM-generated rationales were positively rated by crowdworkers, questions remain regarding the effectiveness of such knowledge grounding. To evaluate whether any fragments of the rationales generated using our proposed approach were grounded on the retrieved knowledge facts, we conducted an experiment. We primarily focus on the corroboration component as there is a higher probability of the knowledge graph containing facts about the correct answer choice.

We measure the existence of knowledgegrounding as follows: consider the retrieved knowledge corresponding to the correct choice j for question q_i in dataset D, \mathcal{G}_{ij} , and the corroboration component of the corresponding LLMgenerated rationale, RC_i . We first measure the BERTScore (Zhang et al., 2019) similarity between a fact $f \in \mathcal{G}_{ij}$, expressed in natural language and a sentence $s \in RC_i$. We then select the fact-sentence pair, (f, s), with the highest BERTScore as a po-

Α	greement = yes (†)		Agreement = no (*) Agreemen			greement = unsure		
Accuracy 66%	Accuracy 90%	Stat. Sig.	Accuracy 66%	Accuracy 90%	Stat. Sig.	Accuracy 66%	Accuracy 90%	Stat. Sig.
$\eta = 6.00$	$\eta = 7.00$		$\eta = 2.00$	$\eta = 3.00$		$\eta = 4.00$	$\eta = 5.00$	
$\mu = 5.89$	$\mu = 6.29$	$\mathbf{p} < 0.01$	$\mu = 2.23$	$\mu = 3.13$	$\mathbf{p} < 0.05$	$\mu = 4.00$	$\mu = 5.25$	p > 0.05
$\sigma = 1.62$	$\sigma = 1.33$		$\sigma = 1.29$	$\sigma = 1.64$		$\sigma = 1.41$	$\sigma = 1.03$	

Table 9: Participant feedback on individual task indicates a more negative impression — rated on a scale between 1 (misleading) to 7 (helpful) — regarding the corresponding rationale. (†) indicates statistical significance with pa < 0.01 and (*) indicates statistical significance with p < 0.05.

Metric	Confid	ence (†)	Reliab	ility (†)	Safe	ty (†)	Satisf	action	Accepta	ability (†)
Accuracy	66%	90%	66%	90%	66%	90%	66%	90%	66%	90%
Median	3.00	4.00	2.00	4.00	3.00	4.00	3.00	5.00	3.00	4.0
Mean	2.91	4.09	1.82	3.64	2.45	3.72	3.45	4.55	3.09	4.27
Std. Dev.	1.14	0.54	0.87	1.03	1.13	0.79	1.44	0.52	1.04	0.65

Table 10: Participant feedback on individual task indicates a more negative impression regarding the corresponding rationale. (†) indicates statistical significance with p < 0.01.

Metric		Agreement %	0
Metric	Overall	Accuracy 66%	Accuracy 90%
Agreement = yes	76.67%	67.27%	86.07%
Agreement = no	20.30%	31.52%	9.09%
Agreement = unsure	3.03%	1.21%	4.85%

Table 11: Participant feedback on individual task indicates a higher disagreement with lower confidence model prediction.

Dataset	Pairwise Max BERTScore	Percentage of Entailment
CSQA	$\mu = 0.5823, \sigma = 0.0650$	80.4%
OBQA	$\mu = 0.5173, \sigma = 0.0803$	38%

Table 12: Degree of knowledge grounding observed in the LLM-generated rationales.

tential candidate for evaluating whether the fact *f* entails the sentence *s* within the rationale. Such entailment is an indicator of whether a fragment of a rationale being grounded on retrieved knowledge facts. Similar approach has been adopted in existing work (Wu et al., 2023) to extract candidate sentences from long documents and evaluate the degree to which the corresponding summary is grounded on the source document. Following their approach, we employ NLI models (Reimers and Gurevych, 2019), *i.e.*, DeBERTa-base model fine-tuned on SNLI (Bowman et al., 2015) and MNLI (Williams et al., 2018), to evaluate entailment. For the BERTScore, we used DeBERTa-Large model (He et al., 2020) fine-tuned on MNLI.

1370

1371

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1388 1389

1390

1391

1393

We measure the knowledge-grounding statistics of the CSQA and OBQA dataset rationales evaluated in the acceptability crowd study in Section 5. As shown in Table 12, on average, at least one fact-sentence pair achieved BERTScore of 0.5823 and 0.5173 for CSQA and OBQA datasets, respectively. While a higher percentage of those pairs were classified as entailment (80.4%) for CSQA, the entailment statistics was a bit lower for OBQA. On reflection, the lower value seems reasonable since we used ConceptNet, a commonsense knowlege graph, as the external source for OBQA, a dataset on elementary science question answering.

The initial observations highlight the promise of knowledge-guided rationalization in ensuring factuality of LLM-generated rationales. However, more in-depth analysis with a stronger metric that takes into account multiple fact-sentence pair candidates across corroboration and refutation components is required to reliably capture the degree of knowledge. Such fine-grained analysis is beyond the scope of our study and can be explored in future.

D.2 A Deeper Dive into LLM vs ECQA

To better understand, we further analyze the crowd 1408 worker feedback based on their preference of ra-1409 tionales. Cases where workers preferred LLM-1410 generated rationales over humans (*i.e.*, the 61.8%1411 cases) - LLM-generated rationales were rated sub-1412 stantially higher than human-written rationales, ex-1413 cept conciseness (see Figure 8.) Even the con-1414 ciseness rating for both types of rationales was 1415 almost the same, with human-written rationales far-1416 ing slightly better. On the other hand, for cases 1417 where workers preferred human-written rationales 1418 over LLMs (*i.e.*, the 38.2% cases) — surprisingly, 1419 apart from conciseness, human-written rationales 1420 were rated significantly higher only on two aspects: 1421 factuality and convincingness. For the rest of the 1422 aspects, the differences between ratings of both 1423 rationale types were marginal. 1424

E Study Interfaces

1425

1394

1395

1396

1398

1399

1400

1401

1402

1403

1404

1405

1406

In this section, we provide screenshots of the important aspects of the three studies.



Figure 8: (a) LLM-generated rationales preferred over human-written (ECQA) rationales. LLM-generated rationales were rated substantially higher than human-written rationales, with the exception of conciseness. (b) ECQA rationales preferred over LLM-generated rationales. Surprisingly, human-written rationales were rated significantly higher only on three aspects: conciseness, factuality and convincingness.

E.1 Faithful Rationalization Interface Details

Both studies were conducted in the Amazon 1429 Mehcanical Turk. We mentioned the worker inclu-1430 sion criteria in Section 3. Each study was launched 1431 in separate batches and were not conducted simul-1432 1433 taneously. Due to the complexity of HITs in each of the studies, we designed the study interfaces 1434 from scratch using HTML and JavaScript. These 1435 interfaces were uploaded in the platform as a new 1436 project to launch the corresponding study. 1437

1428

Figure 9 shows a screen shot of the HIT inter-1438 face of the first study - head-to-head comparison 1439 between LLM-generated and ECQA (crowdworker-1440 written) rationales. The HIT contains a question 1441 and the choices, a selected answer, and two ratio-1442 nales, order of which the are determined at random 1443 on-the-fly. Figure 10 shows a screen shot of the 1444 HIT interface of the acceptability crowd study with 1445 a question and the choices, a selected answer, and 1446 an LLM-generated rationales. For both the studies, 1447 the workers were asked several rating questions de-1448 signed to collect feedback on both coarse-grained and fine-grained aspects of a rationale outlined in Section 4 and Section 5. Workers were asked to rate the rationale(s) using a sliding scale ([1, 7]).

1449

1450

1451

1452

1453

1455

1457

1459

1460

1461

1462

1463

1464

We opted for Likert scale-based rating rather than choice questions to get a more fine-grained 1454 feedback. Given a choice questions, each choice may not exactly capture the participants interpreta-1456 tion of how much a rationale observed the property being evaluated. For example, as shown in Figure 9, 1458 we ask the crowdworker to "rate how understandable each rationale is ". To assist the participants, we suggest how to use the scale - provide interpretation of three points in the scale, i.e., 1 = Notunderstandable, 4 = Somewhat understandable, and 7 =Completely understandable.

Additional quality control measures. Note that 1465 1466 some instances in CSOA have multiple correct or very similar answer choices, due to noise in the 1467 dataset and the fact that the wrong answer choices 1468 were deliberately collected to make the task chal-1469 lenging. To remove this possible confounder, fol-1470 lowing related work (Wiegreffe et al., 2022), in 1471 both the crowd studies we instruct crowdworkers 1472 to treat the selected answer as correct even if they 1473 disagree with it, and then rate the fine-grained as-1474 pects of the rationales. To minimize bias, we ran-1475 domized the order in which rationales were dis-1476 played side-by-side across workers of each HIT. 1477 We also randomized randomized the order of the 1478 rating questions on the fine-grained aspects pre-1479 sented across workers of each HIT. Three different 1480 workers completed each HIT. The workers who par-1481 ticipated in the comparative study were excluded 1482 from the acceptability study. Furthermore, for the 1483 acceptability study, we launched the OBQA dataset 1484 HITs after the conclusion of the CSQA HITs and 1485 excluded workers participating in the CSQA HITs. 1486

Credible Rationalization Interface Details 1487 **E.2**

As shown in Figure 11, participants are first asked 1488 to answer a multiple choice question sampled ran-1489 domly from the CSQA and OBQA datasets. We 1490 ensure the accurate distribution of questions with 1491 correct and incorrect KIT model prediction for 1492 each study condition by grouping questions in each 1493 1494 dataset by prediction accuracy. Once the participant selects an answer, they are shown the KIT 1495 model prediction and the LLM-generated rationale 1496 (Figure 12). At this point, the QA component is 1497 disabled so the the participants cannot change their 1498

options. Finally, participants are provided two fol-1499 low up questions to collect immediate feedback 1500 regarding the task (Figure 13). Finally, participant 1501 conclude the study by completing a survey with 1502 questions adapted from the XAI trust scale (Hoff 1503 and Bashir, 2015) (see Figrue 14.) 1504

<i>Question:</i> What is a child likely to do w	vhile going to play?					
Choices: a) laugh b) sit c) happiness d) being entertained e) walk slowly					
Selected Answer: laugh						
Rationale 1	Rationale 2					
Children are generally joyful and happy while going to play and people who are joyful and happy usually laugh a lot. So a child is likely to laugh while going to play. Happiness comes under laugh and people sit and get entertained in the play, not while going to play. People generally don't walk slow to play.	The answer is laugh because commonsense suggests that children typically laugh when they are going to play. This is because playing is often associated with fun and laughter. While sitting, being entertained, walking slowly, and feeling happiness are all associated with going to play, they are not the primary action that a child is likely to do.					
Rate how <u>convincingly</u> does the rationaliz	zation refute the unselected choices.					
1 = Not convincing, 4 = Somewhat c	onvincing, 7 = Very convinving					
Selected Slider Value: 0	Selected Slider Value: 0					
Not convincing O Rationale 1 Very convincing	Not convincing O Rationale 2 Vary convincing					
Rate how <u>understandab</u>	Rate how <u>understandable</u> each rationale is:					
1 = Not understandable, 4 = Somewhat unders Selected Silder Value: 0	tandable, 7 = Completely understandable Selected Slider Value: 0					
Not understandable O	Not understandable O Completely understandable Rationale 2					

Figure 9: A partial screenshot of the head-to-head comparison interface.

Study Instructions Examples
Question: What is a child likely to do while going to play?
Choices: a) laugh b) sit c) happiness d) being entertained e) walk slowly
Selected Answer: laugh
Rationale
The answer is laugh because commonsense suggests that children typically laugh when they are going to play. This is because playing is often associated with fun and laughter. While sitting, being entertained, walking slowly, and feeling happiness are all associated with going to play, they are not the primary action that a child is likely to do.
Rate how grammatically correct the rationale is:
1 = Not grammatically correct, 4 = Somewhat grammatically correct, 7 = Grammatically correct Selected Slider Value: 0
Not grammatically correct O Completely grammatically correct

Figure 10: A partial screenshot of the acceptability task interface.

	Welcome	HIT Instruction	Quiz	Final Survey	End
Task 1/15					
Instructions (click to ex	(pand/collapse)				
Question: Every evening	ı a child can lool	k into the night sky a	and see that the	moon is	
Choices: a) gone b) brea	aking c) falling d) moving upwards			
Your answer:					
🔵 a) gone 💿 b) breal	king 💿 c) fall	ing 💿 d) moving	upwards		
Submit answer					

Figure 11: For each task, participants are first asked to answer a multiple choice question.

	Welcome		Quiz	Final Survey	End	
	weicome	HIT Instruction	Quiz	Final Survey	Enu	
Task 1/15						
Instructions (click to ex	pand/collapse)					
Question: Every evening	a child can look	k into the night sky a	and see that the	moon is		
Choices: a) gone b) brea	king c) falling d) moving upwards				
Your answer:						
🔵 a) gone 💿 b) break	ting 💿 c) falli	ing 🛛 O d) moving	upwards			
AI prediction: d) moving	upwards					
AI rationale: The answer night sky changes every o elliptical orbit around the	is "moving upw evening. The m Earth.	vards" because the oon is not gone, fall	moon moves in ing, or breaking	an elliptical orbit . It is not stationa	around the Ear ary, but rather is	th, and its position in the constantly moving in an
Show Follow-up						

Figure 12: Once the participant selects an answer, they are shown the KIT model prediction and the LLM-generated rationale.

	Welcome	HIT Instruction	Quiz	Final Survey	End				
Task 1/15									
Instructions (click to ex	pand/collapse)								
Question: Every evening	a child can lool	k into the night sky	and see that the	e moon is					
Choices: a) gone b) brea	aking c) falling d) moving upwards							
Your answer:									
🔵 a) gone 🔷 b) break	<mark>king 🕕 c) falli</mark>	ing 🛛 O d) moving	upwards						
AI prediction: d) moving	upwards								
Al rationale: The answer night sky changes every elliptical orbit around the	r is "moving upw evening. The m Earth.	vards" because the oon is not gone, fal	moon moves in ling, or breaking	an elliptical orbit I. It is not stationa	around the Ear ary, but rather is	th, and its position in the s constantly moving in an			
Show Follow-up									
Do you agree with the A	Al prediction?								
🔾 Yes 📄 No 📄 Un	decided								
How would you charact	erize the AI rat	ionale's role in jus	tifying the AI p	rediction?					
actively misleading	1 2 0	3 4 5	0 6 7	helpful and expl	anatory				
Next Page									

Figure 13: Collecting immediate participant feedback for a task.

Please express to what extent you agree with the following statements based on your experience in the quiz phase.

1. I am confident in the AI rationalizer. I feel that it works well.

Strongly disagree 1 2 3 4 5 Strongly agree

2. I like using the AI rationalizer for understanding decision making process of a AI model.

Strongly disagree 1 2 3 4 5 Strongly agree

3. Overall the AI rationalizer generates highly acceptable rationales of AI predictions.

Strongly disagree 1 2 3 4 5 Strongly agree

4. The AI rationalizer generating the rationales is very reliable. I can count on it to be correct all the time.

Strongly disagree 1 2 3 4 5 Strongly agree

5. This is an attention check. Please select 5 as your response. The AI rationalizer is efficient in that it works very quickly.

Strongly disagree 1 2 3 4 5 Strongly agree

6. The AI rationalizer can perform the task better than a novice human user.

Strongly disagree 1 2 3 4 5 Strongly agree

7. I feel safe that when I rely on the AI rationalizer I will get the right rationales.

Strongly disagree 1 2 3 4 5 Strongly agree

8. The rationalization process of the AI rationalizer is very predictable.

Strongly disagree 1 2 3 4 5 Strongly agree

Next Page

Figure 14: Trust scale-based survey of participant experience.