Reconciling Fidelity and Expressiveness in Large Language Models

Anonymous ACL submission

Abstract

Grounding on external knowledge to generate responses is an effective method to mitigate hallucinations for Large Language Models (LLMs). However, current LLMs struggle to 004 weave knowledge into responses seamlessly while ensuring fidelity like humans, often producing outputs that are either unsupported by external knowledge or overly verbose and unnatural. In this work, to break the trade-off between fidelity and expressiveness, we propose Collaborative Decoding (CoDe), which integrates the output probabilities with and without external knowledge based on their distribution divergence and model confidence to dynamically arouse relevant and reliable expressions from model's internal parameter. Additionally, a knowledge-aware reranking mechanism is de-017 signed to prevent the model from being overly confident in its prior parameter knowledge and from ignoring the given external knowledge. With extensive experiments, our plug-and-play CoDe achieved excellent performance in enhancing fidelity without sacrificing expressiveness on different LLMs and metrics, proving its effectiveness and generality.¹

1 Introduction

027

033

Although large language models (LLMs) have achieved remarkable performance across a variety of tasks in recent studies (Bai et al., 2023; Yang et al., 2023a; Touvron et al., 2023; OpenAI, 2023a,b), they remain prone to hallucination, generating content that appears plausible but is factually incorrect (Ji et al., 2023; Huang et al., 2025). Research suggests that this issue faced by LLMs stems from their knowledge boundaries (Ren et al., 2023), lack of long-tail knowledge (Kandpal et al., 2023), and outdated internal knowledge. These limitations hinder the practical application of LLMs. To mitigate such issues, augmenting LLMs with external knowledge by incorporating it into the model's



Figure 1: Examples exhibits the trade-off between expressiveness and fidelity in LLMs. Higher x-coordinates correspond to higher fidelity, and higher y-coordinates correspond to better expressiveness. Examples (a), (c), and (d) are constrained by the trade-off, whereas our approach break it and generate responses like (b).

input for reference is a promising way, which significantly improved the factuality of the generated content. In particular, Retrieval-Augmented Generation (RAG) paradigm has been widely used.

However, such external-knowledge-augmented LLMs like RAG still face two significant challenges. First, current LLMs could generate content that contradicts or is unsupported by the external knowledge, as illustrated in response (a) and (c) of Figure 1. Second, current LLMs struggle to seamlessly weave external knowledge into response like humans, often resulting in issues such as poor interactivity, dullness, and redundancy (Chen et al., 2023; Yang et al., 2023b). As shown in response (d) of Figure 1, the LLM simply replicates the provided 041

¹The code will be released at GitHub upon publication.

knowledge rather than offering a direct reply to the user's greeting. This approach significantly reduces 057 the interaction's engagement. Numerous methods have been proposed to alleviate the first challenge (Deng et al., 2023; Sun et al., 2023b; Zhang et al., 2024; Liang et al., 2024), but they ignore the second 061 challenge and even make it worse. A user-preferred 062 LLM assistant must generate responses that are grounded in the given knowledge, which we define as fidelity (or faithfulness), and it must also creatively utilize external knowledge to produce natural, diverse, and engaging responses, which 067 we term expressiveness. Unfortunately, Chawla et al. (2024) discovered a trade-off between these two objectives by experimentally masking the dialogue history and the external knowledge in the model's input, revealing the inherent challenges in balancing fidelity and expressiveness. We expanded this investigation on decoding strategies 074 and discovered that there is a similar trade-off in 075 section 3.2. Content generated using deterministic decoding suffers from poor expressiveness, while stochastic decoding results in low fidelity. 079

To break the trade-off between fidelity and expressiveness for LLMs, we propose a novel Collaborative Decoding (CoDe) method. CoDe integrates the output probabilities with and without external knowledge based on their distribution divergence and model confidence to dynamically arouse relevant and high-confidence natural expressions from model's internal parameter. Specifically, we utilize Jensen-Shannon Divergence (JSD) to characterize the distribution differences and combine local confidence (the highest probability) with global uncertainty (distribution entropy) to measure model confidence, facilitating complementary cooperation between two distributions. This approach, which improves expressiveness without relying on the introduction of randomness, effectively avoids hallucinations resulting from stochastic sampling. Furthermore, to prevent the model from being overly confident in its prior parameter knowledge and ignoring retrieval information, we designed a knowledge-aware reranking mechanism for CoDe. This mechanism reranks the top-k suitable tokens based on their awareness of external knowledge (considering both semantic and attention perspectives), thereby ensuring fidelity.

082

086

094

096

100

103

104

105

106

Our contributions are summarized as follows:

• We probe the trade-off between expressiveness and fidelity for external-knowledgeaugmented LLMs on decoding strategies.

• We propose a novel method CoDe, which effectively boosts both fidelity and expressiveness for LLM assistants in knowledgegrounded scenarios, without introducing any additional training, model, or reflection. 107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

• We validated CoDe's generalization and superiority through experiments compared with ten baseline decoding methods on six LLMs, three datasets, and nine automatic metrics.

2 Related Work

2.1 Hallucinations in Text Generation

Hallucination refers to the generation of LLMs appears plausible but is factually incorrect (Zhang et al., 2023c). Existing research has thoroughly investigated hallucination from multiple perspectives, including its origins (McKenna et al., 2023; Dziri et al., 2022b), detection methods (Zhang et al., 2023a; Manakul et al., 2023; Fadeeva et al., 2023), and mitigation strategies (Choi et al., 2023; Chuang et al., 2023; Li et al., 2023a). An effective solution to mitigate hallucination is resorting to the external knowledge, commonly known as Retrieval-Augmented Generation (RAG). Some studies train models to produce fewer hallucinations by constructing training data, such as preference-aligned data or human-annotated data, to obtain groundtruth responses that focus more on fidelity (Liang et al., 2024; Zhang et al., 2024). Some CoT-based works (Zhou et al., 2022; Chae et al., 2023; Yu et al., 2024) externalize implicit knowledge from the backbone LLM to assist generation in the manner of Chain-of-Thought (Wei et al., 2022) or perform self-reflection (Asai et al., 2024). Compared with them, our CoDe serves as a nearly free-lunch for alleviating hallucination, which does not require additional training, model, or prolonged reflection.

2.2 Generation Decoding Strategy

Decoding strategies focus on how to select the next token from the probability distribution of the vocabulary. Popular decoding approaches include greedy decoding, beam search, and top-k sampling. Nucleus (top-p) sampling considers a dynamic number of words that cumulatively reach the probability p (Holtzman et al., 2020). These stochastic strategies promote diversity at the cost of semantic consistency (Basu et al., 2021; Su et al., 2022), and are correlated with hallucination (Dziri et al.,

2022a). Recently, decoding methods based on con-155 trastive probability distributions have attracted con-156 siderable attention in the community. Contrastive 157 decoding (Li et al., 2023b) maximizes the differ-158 ence between expert log-probabilities and amateur log-probabilities to improve fluency and diversity. 160 DoLa (Chuang et al., 2023) contrasts the logits 161 of mature layer and pre-mature layers to mitigate 162 hallucinations in LLMs. CAD amplifies the differ-163 ence between output probabilities with and with-164 out the context document to highlight the external knowledge (Shi et al., 2024). VCD contrasts output 166 distributions derived from original and distorted visual inputs to mitigate the object hallucinations in 168 Large Vision-Language Models (Leng et al., 2023). 169 Similar to these contrastive approaches, CoDe also 170 involves the integration of two distributions. Previ-171 ous methods failed to achieve a win-win situation 172 in terms of fidelity and expressiveness, while CoDe 173 achieves that by complementary collaboration. 174

3 Preliminaries

175

176

177

178

179

181

183

185

189

190

191

193

194

195

3.1 Task Formulation

We consider an LLM parametrized by θ . The input to the LLM consists of a task-specific instruction \mathcal{I} , a dialogue history h across multiple dialogue turns, the user utterance of the current turn u, and relevant external knowledge $\mathbf{k} = (k_1, \ldots, k_m)$ for the current turn, where m represents the number of tokens in \mathbf{k} . For convenience, we define the conversation context $\mathbf{x} = [\mathcal{I}; \mathbf{h}; \mathbf{u}]$. At each time step, the LLM predicts the next token based on the input and the previously generated tokens $\mathbf{y}_{< t}$, yielding logits over the vocabulary:

$$\operatorname{logit}_{\theta}(y_t|\cdot) = \mathcal{LLM}_{\theta}(\boldsymbol{x}, \boldsymbol{k}, \boldsymbol{y}_{< t}).$$
(1)

Next, the probability distribution at step t could be obtained by passing through a softmax layer. Finally, based on the probabilities $p(y_t|\cdot)_{\theta}$, several decoding strategies are employed to select the next token y_t .

$$y \sim p_{\theta}(y_t | \boldsymbol{x}, \boldsymbol{k}, \boldsymbol{y}_{\leq t}) \propto \exp \operatorname{logit}_{\theta}(y_t | \cdot).$$
 (2)

3.2 Pilot Observations and Insights

196There remains considerable potential for im-197provement in expressiveness and fidelity. Pro-198viding external knowledge as a reference for LLM199could reduce the expressiveness of the generated200responses. LLMs often take a shortcut by directly



Figure 2: The trade-off between fidelity and expressiveness of current decoding strategies on Qwen2.5-chat models at different scales. The dashed line indicates the expressiveness score without referring to knowledge.

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

228

229

230

231

232

extracting and copying external knowledge to generate informative responses, as shown in (c) and (d) of Figure 1. Figure 2 indicates that, when provided with external knowledge, the expressiveness score of generated responses declines, suggesting that LLMs tend to prioritize delivering the external knowledge to the user and overlook to maintain coherence, naturalness, and engagingness. There exists knowledge Fidelity Gap Between LLMs and Humans. LLMs may generate content that contradicts the external knowledge due to erroneous reasoning or prioritization of their internal knowledge, as shown in (a) and (c) of Figure 1. Figure 8 indicates that the fidelity of the responses generated by several widely-used open-source LLMs still lags behind that of human-generated content. There exists trade-off Between Expressiveness and Knowledge Fidelity on current decoding strategies. The results in Figure 2 indicate that the content generated using deterministic decoding suffers from poor expressiveness, while that produced through stochastic decoding often lacks fidelity. In addition, smaller models are more affected by decoding strategies. This paper aims to break the trade-off and achieve a win-win situation for fidelity and expressiveness.

4 Approach

In this section, we introduce the Collaborative Decoding (CoDe) method for external-knowledgeaugmented LLM assistant, which consists of two key components, as shown in Figure 3.

4.1 Adaptive Dual-Stream Fusion

As we analyzed in Section 3.2, when provided with external knowledge in the input, the model



Figure 3: An overview of the CoDe method. CoDe consists of two key components: Adaptive Dual-Stream Fusion Module fuses two complementary generation streams adaptively based on model confidence and distribution divergence to integrate relevant and reliable internal knowledge with external knowledge; Knowledge-Aware Reranking Module selects high-fidelity tokens based on semantic and attentive rewards to prevent the model from being overly confident in its prior parameter knowledge and ensure fidelity.

tends to copy knowledge fragments, diminishing
its expressive power. Although stochastic decoding
methods such as top-k (Fan et al., 2018) and nucleus sampling (Holtzman et al., 2020) could alleviate this issue, their inherent probabilistic sampling
nature inevitably compromises fidelity, leading to
severe hallucinations. Therefore, we aim to design a method beyond introducing randomness to
enhance expressiveness while maintaining fidelity.

Distribution Collaboration. Inspired by the design of contrastive decoding (Li et al., 2023b), we propose a novel dual-stream fusion decoding ap-246 proach. In contrast to their work, our approach em-247 phasizes complementary collaboration between the two output distributions, each focusing on differ-249 ent aspects, rather than filtering out errors through 250 contrast. Specifically, CoDe generates two distinct output distributions: one conditioned only on the conversation context x (expressiveness-oriented stream), and the other on both the relevant ex-254 ternal knowledge k and conversation context x255 (faithfulness-oriented stream). Then, a new collaborative probability distribution is computed by fusing two steams to break the trade-off, and is formulated as:

$$p_{CoDe}(y_t) = \operatorname{softmax}[\alpha \operatorname{logit}_{\theta}(y_t | \boldsymbol{x}, \boldsymbol{k}, \boldsymbol{y}_{< t}) + (1 - \alpha) \operatorname{logit}_{\theta}(y_t | \boldsymbol{x}, \boldsymbol{y}_{< t})],$$
(3)

where larger α indicates more weight on the faithfulness-oriented stream. The Equation 3 could also be written as:

$$p_{CoDe} \propto p_{\theta}(y_t | \boldsymbol{x}, \boldsymbol{y}_{< t}) \left(\frac{p_{\theta}(y_t | \boldsymbol{x}, \boldsymbol{k}, \boldsymbol{y}_{< t})}{p_{\theta}(y_t | \boldsymbol{x}, \boldsymbol{y}_{< t})} \right)^{\alpha}.$$
(4)

In this formulation, we model the prior knowledge as $p_{\theta}(y_t | \boldsymbol{x}, \boldsymbol{y}_{< t})$, which predcits the response solely based on the model's internal knowledge. The posterior probability $p_{\theta}(y_t | \boldsymbol{x}, \boldsymbol{k}, \boldsymbol{y}_{< t})$ denotes the probability updating given the observation of external knowledge \boldsymbol{k} . CoDe dynamically adjusts the model's output probability distribution using pointwise mutual information (PMI) between knowledge \boldsymbol{k} and generation y_t , which enhances the probabilities of candidate tokens that have a stronger association with external knowledge in the original distribution.

Adaptive Fusion Weights α . To mitigate the issue of hallucinations caused by low-probability and irrelevant tokens in the prior distribution, we adaptively adjust the parameter α based on model confidence and distribution divergence. When internal knowledge has low relevance to external knowledge or models exhibit high uncertainty, CoDe suppresses the reliance on internal knowledge and emphasizes the role of external knowledge.

$$\alpha = \frac{\delta \cdot \mathcal{C}_t^k}{\mathcal{C}_t^c + \delta \cdot \mathcal{C}_t^k},\tag{5}$$

where C_t^k and C_t^c denotes the confidence of posterior knowledge and the confidence of prior knowl-

265

266

267

268

269

270

271

272

273

274

275

276

263

287

edge, δ denotes the distribution divergence.

Uncertainty can act as a metric to determine when to trust LLMs (Manakul et al., 2023; Huang et al., 2023; Duan et al., 2023). Inspired by Zhang et al. (2023b), we model the factual confidence from two points of view: local confidence and global uncertainty. Local confidence p_{max} is defined as the highest probability among the candidate tokens, while global uncertainty \mathcal{H}_t is defined as the entropy of the entire probability distribution:

$$p_{max} = \max_{y_t \in \mathcal{V}} p(y_t),$$
$$\mathcal{H}_t = -\sum_{y_t \in \mathcal{V}} p(y_t) * \log_2(p(y_t)).$$
(6)

We then synthesize p_{max} and \mathcal{H}_t using the geometric mean function, deriving the confidence score C_t as follows:

$$C_t = \sqrt[2]{\frac{p_{max}}{\mathcal{H}_t + \eta}},\tag{7}$$

where η is a small constant prevents value overflow.

When the prior probability distribution significantly diverges from the posterior probability distribution, it indicates a conflict between internal knowledge and external knowledge. In such cases, we prioritize the external knowledge and reduce the weight of prior distribution. Conversely, if the two distributions are similar, it suggests that they are describing closely relevant content. Under these circumstances, we increase the weight of the prior distribution, leveraging the knowledge acquired during the pre-training phase to generate, thereby enriching the expressiveness of the response content. Specifically, we incorporate the following dynamic parameter δ into the design of α :

$$\delta = \gamma \cdot \exp(\text{JSD}\left(p_c(y_t) \| p_k(y_t)\right)), \quad (8)$$

where $JSD(\cdot, \cdot)$ denotes the Jensen-Shannon Divergence, and γ is a scale factor.

4.2 Knowledge-Aware Reranking

To prevent the model from being overly confident in its prior parameter knowledge and thereby ignoring retrieval information, we further design a knowledge-aware reranking mechanism for CoDe, as shown in the equation below:

$$\hat{p}_{CoDe}(y_t) = \operatorname{topK}\left\{ (1-\beta) \, p_{CoDe}(y_t) + \frac{\beta}{2} \left[\max_{k_i \in \mathbf{k}} \{ \operatorname{sim}(h_{y_t}, h_{k_i}) \} + \max_{k_j \in \mathbf{k}} \{ \operatorname{att}(y_t, k_j) \} \right] \right\},$$
(9)

where larger β values indicate a stronger amplification of fidelity, h denotes the hidden states, $sim(\cdot, \cdot)$ denotes the cosine similarity, and $att(y_t, k_j)$ denotes the attention weight between y_t and k_j after max-pooling for all the layers and attention heads.

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

351

352

353

354

355

357

358

360

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

The knowledge-aware reranking mechanism guarantees fidelity from two respects: (1) semantic reward: the design prefers tokens that exhibit high semantic similarity to tokens in the external knowledge; (2) attentive reward: the design prefers tokens that pay more attention to the knowledge segment. Semantic rewards are calculated as the cosine similarity between the candidate token and all knowledge tokens, while attentive rewards are derived by selecting the highest attention score on all knowledge tokens. As illustrated in Figure 3, when there is a conflict between the model's internal knowledge and external knowledge (LLM indicates that "Jordan played for 14 seasons", while the external knowledge suggests that "he played for 15 seasons"), knowledge-aware reranking mechanism strengthens the awareness of external knowledge, thereby selecting the correct result $(14 \rightarrow 15)$.

The final prediction y_t is obtained by selecting the token that achieves the highest overall fidelity and expressiveness among the top-K candidates:

$$y_t = \arg\max\hat{p}_{CoDe}(y_t). \tag{10}$$

5 Experiments

5.1 Experimental Setup

Datasets. We conducted experiments on three information-seeking dialogue datasets: FAITH-DIAL (Dziri et al., 2022a), HalluDial (Luo et al., 2024), and Wizard of Wikipedia (WoW) (Dinan et al., 2018). These datasets provide dialogue contexts and external knowledge to generate responses, matching the task scenario we aim to investigate. For a detailed description of the datasets, please refer to the Appendix B.

Models. We selected six representative LLMs with varying sizes, architectures, and training data for evaluation, including Llama2-7B-chat (Touvron et al., 2023), Llama-3.1-8B-chat (Grattafiori et al., 2024), Mistral-7B-Instruct-v0.2 (Jiang et al., 2023), Qwen-2.5-3B-chat, Qwen-2.5-7B-chat, Qwen-2.5-14B-chat (Qwen et al., 2025). For implementation details, please refer to the Appendix D.

Baselines. We choose ten decoding methods as the baselines. **Search Methods:** Greedy Decoding (Greedy), Beam Search (Beam) (Boulanger-Lewandowski et al., 2013), Contrastive Search (CS)

289

290

291

298

301

304

306

310

311

312

315

316

317

318

321

323

326

	FAITHDIAL							HalluDial						
Method	Exp	pressive	ness	Fa	aithfulness		Expressiveness			Faithfulness				
	DIV	СОН	CRE	F-Critic	H-Judge	K-BP	DIV	СОН	CRE	F-Critic	H-Judge	K-BP		
Greedy	31.4	57.3	30.0	28.5	86.9	60.9	36.9	64.6	30.1	30.2	87.8	61.2		
Beam	30.8	57.6	25.6	31.3	89.3	<u>64.9</u>	36.1	64.4	24.3	31.5	89.0	65.7		
CS	33.9	55.4	30.0	30.9	83.2	58.7	37.5	<u>64.6</u>	30.2	30.4	88.7	60.9		
FECS	32.8	56.8	28.0	31.6	88.1	63.5	39.4	64.4	30.4	31.0	89.9	64.3		
top-k	36.2	57.2	34.5	21.8	75.3	56.8	<u>40.7</u>	63.8	36.0	21.1	73.0	56.4		
Nucleus	<u>35.6</u>	57.2	<u>34.3</u>	23.4	79.7	57.4	39.9	64.1	<u>34.6</u>	25.3	79.0	57.8		
F-Nucleus	34.1	<u>57.3</u>	32.9	24.3	82.0	58.6	38.5	64.4	32.3	25.7	82.4	59.1		
CD	35.0	55.9	31.3	22.6	76.2	57.0	38.4	62.9	30.5	24.1	78.9	57.3		
DoLa	32.8	56.2	32.3	31.4	87.3	61.2	39.0	64.0	33.6	32.2	89.1	60.4		
CAD	29.2	52.8	21.7	32.1	90.4	67.0	35.4	59.8	22.3	<u>33.6</u>	90.4	<u>67.3</u>		
CoDe	<u>35.6</u>	57.6	29.9	32.4	90.8	67.0	40.9	64.9	29.8	34.3	90.4	67.5		

Table 1: Automatic evaluation results on the FAITHDIAL and HalluDial dataset (Llama2-7B-chat). The best results are highlighted with **bold**. The second-best results are highlighted with underline.

Method	BLEU-2/4	METEOR	ROUGE-L
Greedy	17.0/8.0	20.1	27.1
Beam	17.7/8.7	21.0	28.2
CS	16.4/7.6	18.7	25.7
FECS	18.4/8.8	20.8	28.1
top-k	14.3/6.2	18.2	24.0
Nucleus	14.8/6.4	18.5	24.5
F-Nucleus	15.6/7.0	18.9	25.4
CD	13.6/5.9	17.6	24.5
DoLa	18.1/8.7	20.0	27.7
CAD	18.6/8.9	21.2	27.6
CoDe	19.2/9.5	22.6	29.3

Table 2: Automatic Evaluation results on the FAITH-DIAL dataset (Llama2-7B-chat).

(Su et al., 2022), and FECS (Chen et al., 2023). Stochastic Methods: Top-k Sampling (Fan et al., 2018), Nucleus Sampling (Nulceus) (Holtzman et al., 2020), and Factual-Nucleus Sampling (F-Nucleus) (Lee et al., 2023). Contrastive Methods: Contrastive Decoding (CD) (Li et al., 2023b), DoLa (Chuang et al., 2023), and Context-Aware Decoding (CAD) (Shi et al., 2024). The details of the baseline introduction and hyperparameter settings are found in the Appendix C.

5.2 Experimental Results

5.2.1 Automatic Evaluation

We conducted a comprehensive automated evaluation based on nine metrics from three perspectives:
Faithfulness. To evaluate faithfulness, we adopted the BERT-Precision between the knowledge and generated response (K-BP) following Chen et al. (2023); F-Critic, the average entailment score on FaithCritic (an NLI-based finetuned model) introduced by Dziri et al. (2022a); H-Judge, the ratio



Figure 4: Comparison of knowledge utilization patterns among CoDe, CAD and Nucleus. A more focused distribution towards the bottom-right is preferable.

of samples judged to be faithful by Hallujudge (a specialized LLM-as-a-Judge model) (Luo et al., 2024).

Expressiveness. We considered the model's expressiveness in terms of three aspects: diversity (DIV), context coherence (COH), and the creativity in knowledge utilization (CRE). DIV is calculated as the geometric mean of Distinct-n (n=1, 2, 3, 4) (Li et al., 2016). A high DIV score means that the generation is lexically diverse. Following Su et al. (2022) and Li et al. (2023b), we approximated coherence by cosine similarity between the sentence embeddings of context x and generation y. To calculate CRE, we used the COVERAGE divided by the square root of DENSITY (Grusky et al., 2020). A higher CRE score indicates that the response utilizes more knowledge in a non-extractive manner. More details can be found in Appendix E.

Quality. To assess the overall quality of generated responses, we selected three widely-used metrics based on calculating overlap with the ground-truth: BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and ROUGE (Lin, 2004). For all

Model	Method	Expressiveness		Faithfulness			Quality			
model	1010tillou	DIV	СОН	CRE	F-Critic	H-Judge	K-BP	BLEU-2/4	METEOR	ROUGE-L
Llomo 2.1.9D shot	Regular	34.4	53.4	34.2	46.6	92.2	67.7	21.6/10.4	22.2	31.0
Liama-5.1-6D-Chat	CoDe	35.7	54.5	33.8	50.2	94.0	71.7	21.9/10.8	23.5	30.8
Missiel 7D Laster of 0.2	Regular	34.2	59.5	41.3	21.8	89.8	59.9	15.0/6.7	19.6	25.2
MISUAI-/D-IIISUUCI-VO.2	CoDe	35.4	59.9	38.7	24.3	91.3	6 2.7	15.5/6.9	20.6	25.0
Owen 2.5.2D shot	Regular	37.7	52.4	37.6	38.7	90.5	56.2	18.6/8.5	16.2	25.8
Qwell-2.3-5D-chat	CoDe	39.4	54.4	37.0	42.9	92.9	61.3	20.9/9.8	18.9	27.4
Orwar 2.5.7D shat	Regular	36.8	55.7	40.5	36.0	91.2	61.6	16.8/7.6	18.5	25.5
Qwen-2.5-/B-chat	CoDe	37.7	55.8	40.8	39.6	92.8	64.8	17.6/8.0	20.6	26.4
Owen 2.5.14P abot	Regular	37.8	53.6	39.3	36.6	91.9	65.4	21.7/10.3	21.6	30.1
Qweii-2.3-14D-ciiat	CoDe	38.6	53.6	39.6	36.9	92.6	66.2	22.4/10.5	22.0	30.7

Table 3: More automatic evaluation results compared with regular decoding strategy on the FAITHDIAL dataset. **Regular** denotes widely-used greedy search method in dialogue generation.

of the above metrics, higher is better (\uparrow).

422

423

494

425 426

427

428

429

430

431

432

433 434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

As shown in Table 1 and Table 5, our CoDe surpasses all of the ten baseline decoding methods in terms of three faithfulness metrics. The superior performances of CoDe are consistent among three datasets. In terms of diversity and relevance metrics, our approach achieve top-2 performance. The results of the CRE metric indicate that CoDe exhibits less direct copying of knowledge snippets than other fidelity-enhancing methods such as Beam Search and CAD. We delve deeper into the knowledge utilization patterns, which are shown in Figure 4. We observe that CoDe's density distribution is positioned lower than CAD's, and its coverage is higher compared to sampling methods. This suggests that our approach expresses knowledge in a more natural and diverse manner. Tables 2 and 6 demonstrate that CoDe performs more closely to the ground-truth in traditional metrics, indicating its overall better performance. The results in Table 3 show that CoDe significantly improves fidelity and expressiveness across models of different sizes, architectures, and training corpora. Especially on Llama2-7B-chat and Qwen-2.5-3Bchat, our method outperforms the regular greedy method by +3.9% and +2.4% H-Judge score on the FAITHDIAL dataset, respectively. We also notice that CoDe enables the 3B small LLM to outperform larger models on DIV, COH, F-Critic, and H-Judge, demonstrating its effectiveness.

5.2.2 LLMs-based Evaluation

We further resort to GPT-4, a strong closed-source
LLM, to conduct LLM-as-a-Judge style evaluations (Liu et al., 2023; Chiang and yi Lee, 2023;
Zheng et al., 2023). Typically, we randomly sampled 200 data samples from the test set of FAITHDIAL. We then employed five decoding methods



Figure 5: LLM-based evaluation results on the FAITH-DIAL dataset (Llama2-7B-chat).

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

to generate responses and requested GPT-4 to assign scores from 1 to 5 according to the following criteria: Naturalness (Nat.), Coherence (Coh.), Informativeness (Inf.), Creativity (Cre.), Faithfulness (Fai.), and Factuality (Fac.), strictly adhering to the established rating strategy (Fu et al., 2023). More evaluation details are found in Appendix E.2. From Figure 5, we observe that CoDe effectively breaks the trade-off between expressiveness and fidelity, achieving the best overall performance. In contrast, nucleus sampling and fidelity-enhanced CAD, inevitably lean towards one side of the two dimensions. Compared to the regular greedy search method, our Code demonstrates notable enhancements across nearly all dimensions, aligning with the automated evaluation outcomes in Table 3.

5.2.3 Human Evaluation

Considering the potential biases in automated evaluations and LLM-based evaluations, we also conducted human evaluations. We randomly selected 200 samples from the FaithDial test set for evaluation. Given the dialogue context, related knowledge and the responses generated by CoDe and its baselines, five well-educated annotators were asked to choose the superior response based on three cri-



Figure 6: Human evaluation results on the FAITHDIAL dataset (Llama2-7B-chat). The result is statistically significant with p-value < 0.05, and Kappa (κ) falls between 0.5 and 0.7, suggesting moderate agreement.

Setup		Exp	oressive	ness	Faithfulness			
		DIV	СОН	CRE	F-Critic	H-Judge	K-BP	
Α	CoDe	<u>35.2</u>	<u>57.6</u>	29.9	32.4	90.8	<u>67.0</u>	
В	$-\alpha$	34.9	57.5	<u>32.1</u>	30.1	89.2	64.7	
С	-EOS	34.7	56.8	27.3	32.3	90.8	67.3	
D	-Sem	35.0	57.1	29.6	31.4	88.6	64.1	
Е	-Att	<u>35.2</u>	57.5	30.4	30.9	88.3	63.6	
F	-KAR	35.6	58.0	33.9	29.1	85.5	59.3	

Table 4: Ablation study on the FAITHDIAL dataset.



Figure 7: Hyperparameter study on the FAITHDIAL dataset.

teria: Naturalness, Creativeness, and Faithfulness. Detailed evaluation guidelines can be seen in Figure 9. To measure the agreement among different annotators, we calculate the Fleiss' kappa value (Fleiss, 1971). As shown in Figure 6, CoDe generates significantly more faithful responses compared to all baselines. Evaluators preferred CoDe 1.25x more than greedy search and 5.5x more than CAD when evaluating engagingness. As for naturalness, CoDe is preferred 1.5x more than greedy search and 1.9x more than nucleus sampling.

5.3 Ablation Study

484

485

486

487

488

489

490

491

492

493

494

495

501

In this section, we give detailed ablation studies, 496 including key components, the reward weight β , 497 and the scale factor γ . We first performed the 498 ablation study of the key components to verify 499 the role of each module, including Dynamic Fu-500 sion Weight (- α), Expressiveness-Oriented Stream (-EOS), Knowledge-Aware Reranking (-KAR), semantic reward (-Sem), and attentive reward (-Att). The ablation results for the FAITHDIAL datasets on Llama2-7B-chat are presented in Table 4. The experimental results indicate that every module is essential. In Setup A, the two streams are combined with equal weights, resulting in degraded performance of CoDe on both dimensions. This highlights the importance of appropriately utilizing the model's internal knowledge. When we discard EOS, CoDe encounters a similar situation to other fidelity-improving baselines, where the model's expressiveness decreases. Removing the KAR module leads to a slight improvement in expressiveness but causes an unacceptable decline in fidelity. Setup D and E confirm the necessity of both reward designs. According to the result in Figure 7, setting β to 0.6 and γ to 3 achieves the best performance for CoDe. It also indicates that CoDe is generally robust in the varying settings of hyperparameters.

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

5.4 **Oualitative Examples**

We provide several cases that proves CoDe's strong ability on generating informative and interesting responses. Compared it with the problematic responses produced by the baseline, we observe that CoDe significantly reduces the issues of hallucination and unnatural replies. The cases are shown in Table 9, Table 10, and Table 11, and the corresponding explanation is provided in Appendix A.

Conclusion 6

In this paper, we investigate the trade-off between fidelity and expressiveness for external-knowledgeaugmented LLMs on decoding strategies. Next, we introduce CoDe, a novel plug-and-play decoding method that dynamically integrates relevant and reliable internal knowledge with external knowledge. Extensive experiments conducted on three datasets demonstrate that CoDe effectively breaks the trade-off between fidelity and expressiveness.

Limitations

542

544

545

546

548

552

554

557

558

562

564

566

567

568

571

572

582

583

584 585

587

588

589

591

CoDe effectively boosts both fidelity and expressiveness for LLM assistants in knowledgegrounded scenario. However, we acknowledge that our work still has some limitations.

(1) The performance improvement of CoDe on larger models (14B) is weaker compared to smaller models (3B, 7B). We observe that this is because larger models tend to suffer from overconfidence problems.

(2) As a decoding method, CoDe's applicability is not universal and is primarily suited for externalknowledge-augmented LLM assistants that seek to be engaging and informative. Its suitability diminishes in question-answering contexts where the accuracy and truthfulness of content are prioritized.

(3) Due to time and resource constraints, we have not conducted experiments on larger models (70B, 405B). In the future, we plan to supplement this part of our research.

(4) We did not validate the effectiveness of the CoDe on more specific-designed datasets, such as scenarios where internal knowledge completely conflicts with external knowledge.

Ethics Statement

The benchmark datasets we utilized in our experiments-WoW (Dinan et al., 2018), FAITH-DIAL (Dziri et al., 2022a), and HalluDial (Luo et al., 2024)-are all highly respected, open-source datasets collected by crowdsourced workers. They were compiled with strict adherence to user privacy protection protocols, ensuring the exclusion of any personal information. Moreover, our proposed approach is conscientiously designed to uphold ethical standards and promote societal fairness, guaranteeing that no bias is introduced. For the human evaluation component of our study, all participants were volunteers who received comprehensive information about the purpose of the research, ensuring informed consent. Furthermore, all participants were fairly and reasonably compensated for their contributions. We informed the annotators that the data is to be used solely for academic research purposes.

CoDe does not verify the factual accuracy of the provided reference, so there is a risk of potential misuse. Specifically, if the knowledge source contains non-factual information, CoDe is likely to generate misinformation as well. We strongly advise users to carefully verify the knowledge source before using it. We use AI writing in our work, but only for polishing articles to enhance their readability.

References

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-RAG: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Sourya Basu, Govardana Sachithanandam Ramachandran, Nitish Shirish Keskar, and Lav R. Varshney. 2021. Mirostat: a neural text decoding algorithm that directly controls perplexity. In *International Conference on Learning Representations*.
- Nicolas Boulanger-Lewandowski, Yoshua Bengio, and Pascal Vincent. 2013. Audio chord recognition with recurrent neural networks. In *International Society for Music Information Retrieval Conference*.
- Hyungjoo Chae, Yongho Song, Kai Ong, Taeyoon Kwon, Minjin Kim, Youngjae Yu, Dongha Lee, Dongyeop Kang, and Jinyoung Yeo. 2023. Dialogue chain-of-thought distillation for commonsense-aware conversational agents. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5606–5632, Singapore. Association for Computational Linguistics.
- Kushal Chawla, Hannah Rashkin, Gaurav Singh Tomar, and David Reitter. 2024. Investigating content planning for navigating trade-offs in knowledge-grounded dialogue. In Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2316–2335, St. Julian's, Malta. Association for Computational Linguistics.

635

636

637

638

639

640

641

642

643

644

645

646

592 593

733

734

735

736

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

757

758

759

760

761

762

Wei-Lin Chen, Cheng-Kuang Wu, Hsin-Hsi Chen, and Chung-Chi Chen. 2023. Fidelity-enriched contrastive search: Reconciling the faithfulness-diversity trade-off in text generation. In *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing, pages 843–851, Singapore. Association for Computational Linguistics.

647

671

672

673

674

675

676

677

679

681

684

686

690

694

- Cheng-Han Chiang and Hung yi Lee. 2023. A closer look into automatic evaluation using large language models.
- Sehyun Choi, Tianqing Fang, Zhaowei Wang, and Yangqiu Song. 2023. Kcts: Knowledge-constrained tree search decoding with token-level hallucination detection.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. 2023. Dola: Decoding by contrasting layers improves factuality in large language models.
- Yifan Deng, Xingsheng Zhang, Heyan Huang, and Yue Hu. 2023. Towards faithful dialogues via focus learning. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4554–4566, Toronto, Canada. Association for Computational Linguistics.
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2018. Wizard of wikipedia: Knowledge-powered conversational agents. *CoRR*, abs/1811.01241.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. 2023. Shifting attention to relevance: Towards the uncertainty estimation of large language models.
- Nouha Dziri, Ehsan Kamalloo, Sivan Milton, Osmar Zaiane, Mo Yu, Edoardo M. Ponti, and Siva Reddy. 2022a. FaithDial: A faithful benchmark for information-seeking dialogue. *Transactions of the Association for Computational Linguistics*, 10:1473– 1490.
- Nouha Dziri, Sivan Milton, Mo Yu, Osmar Zaiane, and Siva Reddy. 2022b. On the origin of hallucinations in conversational models: Is it the datasets or the models?
- Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, Timothy Baldwin, and Artem Shelmanov. 2023. Lm-polygraph: Uncertainty estimation for language models.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 889–898, Melbourne, Australia. Association for Computational Linguistics.

- Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76:378–382.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. *CoRR*, abs/2302.04166.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,

Praveen Krishnan, Punit Singh Koura, Puxin Xu, 763 Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, 766 Robert Stojnic, Roberta Raileanu, Rohan Maheswari, 767 Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seo-771 hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye 772 Wan, Shruti Bhosale, Shun Zhang, Simon Van-773 denhende, Soumya Batra, Spencer Whitman, Sten 774 Sootla, Stephane Collot, Suchin Gururangan, Syd-775 ney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal 778 Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-784 feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing 788 Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, 790 Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, 792 Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew 795 Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-797 dani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, 799 Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-807 Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David 810 Xu, Davide Testuggine, Delia David, Devi Parikh, 811 Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, 812 Elaine Montgomery, Eleonora Presani, Emily Hahn, 813 814 Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-815 ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat 816 817 Ozgenel, Francesco Caggioni, Frank Kanayet, Frank 818 Seide, Gabriela Medina Florez, Gabriella Schwarz, 819 Gada Badeer, Georgia Swee, Gil Halpern, Grant 820 Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun 822 Habeeb, Harrison Rudolph, Helen Suk, Henry As-823 824 pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim 825 Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,

Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuvigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models.

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

993

996

943

- Max Grusky, Mor Naaman, and Yoav Artzi. 2020. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies.
 - Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration.
 - Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Trans. Inf. Syst.*, 43(2).
 - Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma.
 2023. Look before you leap: An exploratory study of uncertainty measurement for large language models.

901

903

904

905

906

908

909

910

911

912

913

914

915

916

917

918

921

927

929

935

936

937

938

941

- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.
- Byeongchang Kim, Jaewoo Ahn, and Gunhee Kim. 2020. Sequential latent knowledge selection for knowledge-grounded dialogue.
- Nayeon Lee, Wei Ping, Peng Xu, Mostofa Patwary, Pascale Fung, Mohammad Shoeybi, and Bryan Catanzaro. 2023. Factuality enhanced language models for open-ended text generation.
- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Li Bing. 2023. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 13872–13882.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 110–119, San Diego, California. Association for Computational Linguistics.

- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023a. Inferencetime intervention: Eliciting truthful answers from a language model.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023b. Contrastive decoding: Open-ended text generation as optimization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12286–12312, Toronto, Canada. Association for Computational Linguistics.
- Yuxin Liang, Zhuoyang Song, Hao Wang, and Jiaxing Zhang. 2024. Learning to trust your feelings: Leveraging self-awareness in LLMs for hallucination mitigation. In Proceedings of the 3rd Workshop on Knowledge Augmented Methods for NLP, pages 44–58, Bangkok, Thailand. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Shilei Liu, Xiaofeng Zhao, Bochao Li, Feiliang Ren, Longhui Zhang, and Shujuan Yin. 2021. A three-stage learning framework for low-resource knowledge-grounded dialogue generation.
- 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. *CoRR*, abs/2303.16634.
- Wen Luo, Tianshu Shen, Wei Li, Guangyue Peng, Richeng Xuan, Houfeng Wang, and Xi Yang. 2024. Halludial: A large-scale benchmark for automatic dialogue-level hallucination evaluation.
- Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models.
- Nick McKenna, Tianyi Li, Liang Cheng, Mohammad Javad Hosseini, Mark Johnson, and Mark Steedman. 2023. Sources of hallucination by large language models on inference tasks.
- Chuan Meng, Pengjie Ren, Zhumin Chen, Zhaochun Ren, Tengxiao Xi, and Maarten de Rijke. 2021. Initiative-aware self-supervised learning for knowledge-grounded conversations. In *Proceedings* of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '21, page 522–532, New York, NY, USA. Association for Computing Machinery.

OpenAI.	2023a.	Chatgpt.	994
https://op	enai.com/blog/chatgpt/.		995

OpenAI. 2023b. Gpt-4 technical report.

997

- 1013
- 1014 1015 1016
- 1018
- 1019
- 1024
- 1025
- 1026 1027
- 1028 1029
- 1030 1031
- 1032 1033

1034 1035 1036

1038

1039

- 1042 1043
- 1044 1045

1046 1047

1049 1050 1051

1054

Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. 2021. Understanding factuality in abstractive summarization with FRANK: A benchmark for factuality metrics. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4812-4829, Online. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 311-318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Daviheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 technical report.

Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, J. Liu, Hao Tian, Huaqin Wu, Ji rong Wen, and Haifeng Wang. 2023. Investigating the factual knowledge boundary of large language models with retrieval augmentation. In International Conference on Computational Linguistics.

Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemover, and Wen-tau Yih. 2024. Trusting your evidence: Hallucinate less with contextaware decoding. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers), pages 783-791, Mexico City, Mexico. Association for Computational Linguistics.

- Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. 2022. A contrastive framework for neural text generation.
- Weiwei Sun, Pengjie Ren, and Zhaochun Ren. 2023a. Generative knowledge selection for knowledgegrounded dialogues. In Findings of the Association for Computational Linguistics: EACL 2023, pages 2077-2088, Dubrovnik, Croatia. Association for Computational Linguistics.

Weiwei Sun, Zhengliang Shi, Shen Gao, Pengjie Ren, Maarten de Rijke, and Zhaochun Ren. 2023b. Contrastive learning reduces hallucination in conversations. In Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium

on Educational Advances in Artificial Intelligence, AAAI'23/IAAI'23/EAAI'23. AAAI Press.

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

- Ilva Sutskever, Oriol Vinvals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2, NIPS'14, page 3104–3112, Cambridge, MA, USA. MIT Press.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38-45, Online. Association for Computational Linguistics.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. 2023. Sheared llama: Accelerating language model pre-training via structured pruning. arXiv preprint arXiv:2310.06694.
- Lin Xu, Qixian Zhou, Jinlan Fu, Min-Yen Kan, and 1109 See-Kiong Ng. 2022. CorefDiffs: Co-referential and 1110 differential knowledge flow in document grounded 1111 conversations. In Proceedings of the 29th Inter-1112 national Conference on Computational Linguistics, 1113

- 1114 1115
- 1116 1117 1118
- 1119 1120
- 1121 1122 1123
- 1124 1125
- 1126 1127
- 1128 1129 1130
- 1131
- 1132 1133 1134
- 1135 1136 1137
- 1138 1139
- 1140 1141 1142
- 1143 1144
- 1145 1146
- 1148 1149 1150

- 1151 1152 1153
- 1154 1155 1156
- 1157 1158
- 1159 1160
- 1161 1162
- 1163 1164
- 1165 1166 1167
- 1168 1169

1170 1171 1172

- pages 471-484, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, JunTao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. 2023a. Baichuan 2: Open large-scale language models.
 - Chenxu Yang, Zheng Lin, Jiangnan Li, Fandong Meng, Weiping Wang, Lanrui Wang, and Jie Zhou. 2022. TAKE: Topic-shift aware knowledge sElection for dialogue generation. In Proceedings of the 29th International Conference on Computational Linguistics, pages 253-265, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Chenxu Yang, Zheng Lin, Lanrui Wang, Chong Tian, Liang Pang, Jiangnan Li, Qirong Ho, Yanan Cao, and Weiping Wang. 2023b. Multi-level adaptive contrastive learning for knowledge internalization in dialogue generation. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 8002-8015, Singapore. Association for Computational Linguistics.
- Jiong Yu, Sixing Wu, Jiahao Chen, and Wei Zhou. 2024. LLMs as collaborator: Demands-guided collaborative retrieval-augmented generation for commonsense knowledge-grounded open-domain dialogue systems. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 13586-13612, Miami, Florida, USA. Association for Computational Linguistics.
- Haolan Zhan, Lei Shen, Hongshen Chen, and Hainan Zhang. 2021. CoLV: A collaborative latent variable model for knowledge-grounded dialogue generation. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 2250–2261, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Chenghu Zhou, Xinbing Wang, and Luoyi Fu. 2023a. Enhancing uncertaintybased hallucination detection with stronger focus.
- Tianhang Zhang, Lin Qiu, Qipeng Guo, Cheng Deng, Yue Zhang, Zheng Zhang, Chenghu Zhou, Xinbing Wang, and Luoyi Fu. 2023b. Enhancing uncertaintybased hallucination detection with stronger focus. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 915–932, Singapore. Association for Computational Linguistics.

Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, 1173 Lifeng Jin, Linfeng Song, Haitao Mi, and Helen 1174 Meng. 2024. Self-alignment for factuality: Miti-1175 gating hallucinations in LLMs via self-evaluation. In 1176 Proceedings of the 62nd Annual Meeting of the As-1177 sociation for Computational Linguistics (Volume 1: 1178 Long Papers), pages 1946–1965, Bangkok, Thailand. 1179 Association for Computational Linguistics. 1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023c. Siren's song in the ai ocean: A survey on hallucination in large language models.
- Xueliang Zhao, Wei Wu, Chongyang Tao, Can Xu, Dongyan Zhao, and Rui Yan. 2020a. Low-resource knowledge-grounded dialogue generation.
- Xueliang Zhao, Wei Wu, Can Xu, Chongyang Tao, Dongyan Zhao, and Rui Yan. 2020b. Knowledgegrounded dialogue generation with pre-trained language models.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena.
- Wen Zheng, Natasa Milic-Frayling, and Ke Zhou. 2021. Knowledge-grounded dialogue generation with termlevel de-noising. In Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pages 2972–2983, Online. Association for Computational Linguistics.
- Pei Zhou, Karthik Gopalakrishnan, Behnam Hedayatnia, Seokhwan Kim, Jay Pujara, Xiang Ren, Yang Liu, and Dilek Hakkani-Tur. 2022. Think before you speak: Explicitly generating implicit commonsense knowledge for response generation. In *Proceedings* of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1237-1252, Dublin, Ireland. Association for Computational Linguistics.

Case Study Α

In the example shown in Table 9, Greedy Search 1215 mistakenly treated Madonna's cover song as her 1216 original work, leading to a factual mistake. Greedy 1217 Search, Nucleus Sampling, and DoLa all confused 1218 song titles with album names, and generated many 1219 songs that cannot be verified by external knowledge. Although DoLa indeed activated the model's 1221 internal knowledge to produce numerous accurate 1222 song titles, it kept listing them endlessly, result-1223 ing in an excessively long response. CAD failed 1224 to respond to the user's query and simply copied 1225

		WoW-Seen							WoW-Unseen						
Method	Ex	pressive	ness	Faithfulness			Expressiveness			Faithfulness					
	DIV	СОН	CRE	F-Critic	H-Judge	K-BP	DIV	СОН	CRE	F-Critic	H-Judge	K-BP			
Greedy	31.5	58.3	28.7	19.5	84.2	58.1	22.4	58.2	28.8	17.6	86.7	58.6			
Beam	30.7	58.6	24.1	22.1	86.1	62.4	21.4	58.9	24.0	<u>22.9</u>	87.6	62.8			
CS	32.0	58.4	29.0	21.0	84.2	57.6	23.0	57.5	28.6	18.9	85.9	58.4			
FECS	33.4	57.8	28.1	23.5	86.3	61.9	23.3	57.2	28.6	<u>22.9</u>	<u>88.6</u>	61.9			
top-k	35.8	58.2	33.4	15.0	66.5	53.9	<u>29.8</u>	57.7	33.6	14.2	70.1	54.5			
Nucleus	34.9	58.4	<u>33.1</u>	16.3	73.6	54.9	28.2	57.8	32.8	15.9	76.8	55.5			
F-Nucleus	33.9	<u>58.7</u>	31.3	16.7	77.6	56.1	26.6	58.4	31.3	17.5	80.0	56.6			
CD	34.7	57.1	32.9	16.3	68.9	57.0	34.9	56.9	<u>33.0</u>	13.7	73.8	57.5			
DoLa	33.1	58.0	30.4	22.9	84.2	57.9	23.9	57.5	31.3	21.6	85.9	58.4			
CAD	28.1	53.9	18.6	26.4	<u>86.7</u>	65.6	20.7	53.6	21.8	22.6	87.7	<u>64.8</u>			
CoDe	<u>35.2</u>	58.9	27.7	26.8	88.0	<u>65.0</u>	30.1	<u>58.6</u>	28.1	25.6	89.7	65.2			

Table 5: Automatic evaluation results on the WoW dataset (Llama2-7B-chat).

Method	BLEU-2/4	METEOR	RROUGE-L
Greedy	16.5/7.8	20.2	27.6
Beam	16.6/8.3	22.1	28.6
CS	16.4/7.8	20.0	27.5
FECS	18.3/8.8	20.9	28.7
top-k	13.4/5.7	17.7	23.7
Nucleus	14.2/6.3	18.5	24.6
F-Nucleus	14.9/6.5	19.2	25.4
CD	12.4/5.6	16.8	23.9
DoLa	17.1/8.1	19.6	27.4
CAD	18.5/8.9	22.2	28.0
CoDe	18.5/9.1	22.5	28.9

Table 6: Automatic Evaluation results on the HalluDial dataset (Llama2-7B-chat).

external information into the reply, showing weak expression capabilities. By comparison, CoDe not only interacted with the user, but also correctly utilized external knowledge.

B Datasets

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

WoW is collected based on Wikipedia, with one crowd-sourcer acts as a knowledgeable wizard and the other plays the role of an inquisitive apprentice. The objective is to generate responses based on given knowledge snippets, taken from Wikipedia, that are pertinent to the conversation topic. The ground-truth responses in the dataset are annotated by humans based on the best knowledge they selected. We evaluated all the decoding methods on both the test seen and unseen set. The test seen set includes 4,336 samples where the topics were seen in the training set, while the unseen test set includes 4,370 samples where the topics were not seen in the training set (Dinan et al., 2018).
FAITHDIAL is a benchmark for hallucination-free



Figure 8: Pilot experiment.

dialogues, which optimizes the responses in the WoW dataset to be more faithful to knowledge. Subjective and hallucinated information present in the wizard's utterance of WoW data are edited into utterances faithful to the given knowledge in this dataset. We evaluated all the decoding methods on its test set, which contains 3,539 samples (Dziri et al., 2022a). 1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1261

1262

1263

1265

1266

1267

1268

1269

HalluDial is the first comprehensive large-scale benchmark for dialogue-level hallucination evaluation. It is derived from an information-seeking dialogue dataset and covers factuality and faithfulness hallucinations. The benchmark includes 4,094 dialogues with a total of 146,856 samples. We selected 3,000 samples from its 18,357 non-hallucinatory samples as the test set in our experiments.

C Baselines

Beam search selected the most probable k tokens from the probability distribution at each step to expand the search space (Boulanger-Lewandowski et al., 2013; Sutskever et al., 2014). We set the beam size to 4 in our experiment. **CS:** Su et al. (2022) penalized previously generated

tokens to overcome degeneration and enhance

1270 content diversity. We set $k|\alpha = 4|0.6$.

1278

1279

1282

1283

1284

1295

1296

1298

1299

1300

1301

1302

1304

1305

1306

1308

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1271**FECS:** Chen et al. (2023) extended Contrastive1272Search by integrating a faithfulness term that1273encourages factuality. We set $k|\alpha|\beta = 4|0.3|0.3$.1274**Top-k Sampling** introduced randomness into the1275generation process by selecting from the top-k1276most likely tokens (Fan et al., 2018). We set1277k = 50 in our experiment.

Nucleus sampling considered a dynamic number of words that cumulatively reach the probability p (Holtzman et al., 2020). We set p = 0.9 in our experiment.

F-Nucleus: Lee et al. (2023) modified Nucleus Sampling by adapting the randomness dynamically to improve the factuality of generation. We set $p|\lambda|\omega = 0.9|0.9|0.7$ in our experiment.

1286**CD:** Li et al. (2023b) maximized the difference1287between expert log-probabilities and amateur1288log-probabilities to improve fluency and diversity.1289We use Llama2-7B-chat as the expert model and1290Sheared-LLaMA-2.7B-ShareGPT as the amateur1291model (Xia et al., 2023). We set the amateur1292temperature τ to 1.0. We select the generated1293tokens using a greedy search on the contrasted1294distributions.

DoLa: Chuang et al. (2023) amplified the factual knowledge in LLM by contrasting the logits from different layers to enhance factuality. We set the *dola_layers* hyperparameter to 'high' in our experiments. We select the generated tokens using a greedy search on the contrasted distributions.

CAD: Chuang et al. (2023) amplified the difference between output probabilities with and without the context document to highlight the external knowledge. We set the hyperparameter α to 1.0. We select the generated tokens using a greedy search on the contrasted distributions.

D Implementation Details

We conducted experiments by utilizing the opensource Hugging Face transformers (Wolf et al., 2020). All experiments are conducted with fewshot prompting (three shots). The three demonstrations are manually selected from the FAITH-DIAL dataset, and they are used consistently for all methods during evaluation. We conducted three experiments for all methods, using a different set of samples in each experiment, and finally took the average of the results to eliminate the impact of randomness. We exhibits a set of demonstrations with task instructions in the Appendix G. As 1320 our focus is on the generation process following 1321 the acquisition of retrieval knowledge, we assume 1322 that the knowledge provided to the model is the 1323 most appropriate. For our CoDe and all baselines, 1324 we directly use manually annotated golden knowl-1325 edge from the three datasets as input in the exper-1326 iments. For the hyperparameters in CoDe, we set 1327 $k|\beta|\gamma = 4|0.6|3$. For other decoding hyperparame-1328 ters, we set them to be the same for all methods. We 1329 set the min new tokens to 5, and the batch size 1330 to 1. 1331

Faithfulness					
BERT-Precision	FaithCritic	Hallujudge			
Expressiveness					
Diversity	Coherence	Creativeness			
Quality					
BLEU	METEOR	ROUGE-L			

Table 7: Automatic Evaluation metircs.

1332

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

E Evaluation Metrics

E.1 Automatic Evaluation Metrics

E.1.1 Faithfulness

To evaluate faithfulness, we adopted three faithfulness evaluation metrics, which has been demonstrated to achieve high correlations with human judgment.

K-BP. We calculated BERT-Precision (Pagnoni et al., 2021) between the external knowledge and generated response (**K-BP**) following Shi et al. (2024) to measure the consistency from the perspective of semantic similarity.

F-Critic. FaithCritic is a faithfulness discrimination model fine-tuned on the **FAITHCRITIC** dataset, which is initialized with RoBERTa-Large. This model outputs the probability of positive and negative labels in the form of a binary-classification Natural Language Inference (NLI) task, where responses with subjective and hallucinatory information are predicted as negative labels. **F-Critic.** is the average entailment score on the FaithCritic model.

H-Judge. H-Judge is the ratio of samples judged to be faithful by the Hallujudge model (Luo et al., 2024). Since the authors did not release the weights of Hallujudge, we trained the Hallujudge model on the HalluDial dataset using Meta-Llama-3-8B with the hyperparameters specified in the paper.

- 1360
- 1362
- 1364
- 1365

1367

1368

1369

1371

1376

1378

1379

1380

1381

1383

1384

1385

1386

1387

1388

1389

1390

1392

1393

1394

E.1.2 Expressiveness

We considered the model's expressiveness in terms 1361 of three aspects: diversity, context coherence, and the creativity in knowledge utilization.

DIV. Diversity (**DIV**) is calculated as the geometric mean of Distinct-n (n=1, 2, 3, 4) (Li et al., 2016):

> **DIV** = $\sqrt[4]{\prod_{n=1}^{4}}$ Distinct-n. (11)

COH. Following Su et al. (2022) and Li et al. (2023b), we approximated coherence by cosine similarity between the sentence embeddings of context x and generation y:

$$\mathbf{COH} = \frac{\mathbf{EMB}(x) \cdot \mathbf{EMB}(y)}{||\mathbf{EMB}(x)|| \cdot ||\mathbf{EMB}(y)||}, \qquad (12)$$

where $EMB(\cdot)$ is the SimCSE sentence embedding 1372 1373 (Gao et al., 2021).

> **CRE.** To calculated **CRE**, we use the COVERAGE divided by the square root of DENSITY (Grusky et al., 2020) as follows:

$$\mathbf{CRE} = \frac{\text{Coverage}}{\sqrt[2]{\text{Density}}},\tag{13}$$

$$\operatorname{Coverage}(k, y) = \frac{1}{|y|} \sum_{f \in \mathcal{F}(k, y)} |f|, \qquad (14)$$

$$\text{Density}(k,y) = \frac{1}{|y|} \sum_{f \in \mathcal{F}(k,y)} |f|^2, \quad (15)$$

where $\mathcal{F}(k, y)$ is the set of shared sequences of tokens in knowledge k and response y. A higher Coverage score indicates more knowledge are integrated into the response, while a lower Density score indicates the knowledge are weaved into response naturally and creatively. To unify the measurement of Coverage and Density, we performed a square root operation on Density.

E.1.3 Quality

To assess the overall quality of generated responses, we selected three widely-used metrics based on calculating overlap with the ground-truth: BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and **ROUGE** (Lin, 2004).

E.2 LLM-based Evaluation Metrics

We perform LLM-based evaluation according to 1397 the following criteria: Naturalness (Nat.), Co-1398 herence (Coh.), Informativeness (Inf.), Creativity 1399

(Cre.), Faithfulness (Fai.), and Factuality (Fac.). 1400 High naturalness refers to the generated content 1401 being realistic, engaging, and interactive, capable 1402 of encouraging users to engage in more rounds of 1403 conversation. High coherence means the generated 1404 content is related to the context and maintains a 1405 smooth flow. High informativeness indicates that 1406 the generated content is rich in information and can 1407 help users acquire new knowledge. High creativity 1408 refers to the model's diverse utilization of external 1409 knowledge, rather than mechanically extracting and 1410 directly outputting it. High faithfulness indicates 1411 that the generated content does not contain informa-1412 tion that directly contradicts the given knowledge, 1413 or cannot be verified from the provided knowledge. 1414 High factuality indicates that the generated content 1415 does not conflict with established world knowledge. 1416 To strictly differentiate between faithfulness and 1417 factuality, we included formal definitions for both 1418 in Appendix H. Table 8 shows the prompts we used 1419 for LLM-based evaluation (LLM-as-a-Judge). Our 1420 evaluation was conducted using gpt-4o-2024-08-1421 06. 1422

E.3 Human Evaluation Metrics

Given the dialogue context, related knowledge and the responses generated by CoDe and its baselines, five well-educated annotators were asked to choose the superior response based on three criteria: Naturalness (Nat.), Engagingness (Eng.), and Faithfulness (Fai.). Detailed evaluation guidelines can be seen in Figure 9.

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

F **Additional Related Work**

F.1 Knowledge-Grounded Dialog Generation

Knowledge-grounded dialogue generation aims to 1433 alleviate dull and unfaithful responses by injecting 1434 external knowledge into input of dialogue mod-1435 els and it consists of two subtasks: knowledge 1436 selection and response generation. The hot spot 1437 of early research is mainly concentrated on how 1438 to improve the performance of knowledge selec-1439 tion (Sun et al., 2023a; Xu et al., 2022; Zhan et al., 1440 2021; Yang et al., 2022; Meng et al., 2021; Kim 1441 et al., 2020). With the remarkable leap in the ca-1442 pabilities of generative models, the research focus 1443 gradually shifts to the response generation subtask 1444 (Zhao et al., 2020a; Liu et al., 2021; Zhao et al., 1445 2020b; Zheng et al., 2021). Ideally, a brilliant robot 1446 should generate informative and truthful responses 1447

1448	while maintaining the naturalistic phrasing and ex-	1983 t
1449	cellent interactivity (Dziri et al., 2022a).	Exam
		User's
1450	G Details of the Generation Prompts	sound
1451	Our instruction template and demonstrations for	aboui ###kn
1452	prompting Large Language Models to generate re-	nated
1453	sponse during evalutation are as follows:	Orlea
1454	As an assistant, your task is to engage in a chit-	20th c
1455	chat conversation with user. You will be provided	and ro
1456	the dialogue history and a piece of related knowl-	Vour k
1457	edge, and your task is to utilize the knowledge to	like ia
1458	continue the conversation. Your English response	comm
1459	should be informative but no more than 50 words.	20th c
1460	coherent with the dialogue context and faithful to	20m c Now c
1461	the given knowledge	Now C
1462	You SHOULD refer to the following examples:	<i>###1</i>
1/63	Frample 1.	###KIL
1/6/	User's utterance: Shower	10UF K
1404	###howledge###: Ancient people heads to repro	п
1400	$\pi\pi\pi\kappa$ how ledge $\pi\pi\pi$. Ancient people began to repro-	пг
1400	unce these natural phenomena by pouring jugs of	We for
1407	water, often very cola, over themselves after wash-	explai
1468	Ing.	Defini
1469	Your knowledge-grounded response to the user: I	and ex
1470	love taking snowers! I do not know now I would	n we
14/1	live without the modern showers. In ancient times	to the
1472	people would pour jugs of cold water over them-	the fol
1473	selves and consider that a shower.	ine joi
1474	Example 2:	• =
1475	User's utterance: Is rock and roll still popular to-	-
1476	day?	n
1477	Your response: It's hard to say. However, radio	ti
1478	stations have much success playing classic rock	
1479	and roll, which is a sub genre that usually has one	g
1480	or two electric guitars, a double bass or string bass	Defini
1481	or electric bass guitar, and a drum kit.	we sa
1482	User's utterance: I used to listen to the rock band	follow
1483	Rolling Stones. Are they still around today?	jouon
1484	Your response: They are! Even though they were	• =
1485	formed in 1962 and have had a long list of line-up	
1486	changes, they're still around today, with Mick Jag-	S
1487	ger still leading the band.	5
1488	User's utterance: Wow, that is a long time to be	Theor
1489	playing music. I wonder if any other bands have	entailt
1490	been around that long.	The
1491	###knowledge###: Red Hot Chili Peppers are an	faithfi
1492	American funk rock band formed in Los Angeles in	Verse 1
1493	1983.	rem in
1494	Your knowledge-grounded response to the user: It	10111 18
1495	all depends! You have bands like the Red Hot Chili	тт
1496	Peppers who, although have not reached the popu-	1 11
1497	larity of the Rolling Stones, have been around since	Theor

1983 themselves.	1498
Example 3:	1499
User's utterance: jazz music is a very interesting	1500
sound and interesting genre. what can you tell me	1501
about it?	1502
###knowledge###: Jazz is a music genre that origi-	1503
nated in the African-American communities of New	1504
Orleans, United States, in the late 19th and early	1505
20th centuries, and developed from roots in blues	1506
and ragtime.	1507
Your knowledge-grounded response to the user: I	1508
like jazz too. It originated in the African-American	1509
communities of New Orleans in the late 19th and	1510
20th centuries.	1511
Now complete the following dialogue:	1512
User's utterance: [Dialogue Context]	1513
###knowledge###: [External Knowledge]	1514
Your knowledge-grounded response to the user:	1515
II Faithfulness and Fastuality	
H Faithuiness and Factuality	1516
We formally define faithfulness and factuality and	1517
explain their relationship to facilitate the research:	1518
Definition 1 Faithfulness(\mathcal{F}): Given a response y,	1519
and external knowledge $\mathcal{K} = (k_1, \dots, k_i)$ at turn	1520
n, we say that the response y is faithful with respect	1521
to the external knowledge $(\mathcal{F}(\mathcal{K}, y))$ if and only if	1522
the following condition holds:	1523
• $\exists \Gamma$ such that $\Gamma \models y$, where Γ is a non-empty	1524
subset of \mathcal{K} and \models denotes semantic entail-	1525
ment. In other words, there is no interpreta-	1526
tion $\mathcal I$ such that all members of Γ are true and	1527
y is false (Dziri et al., 2022a).	1528
Definition 2 Eactuality(\mathcal{T}): Cince a response of	4500
we say that u is factual $(T(u))$ if and only if the	1529
following condition holds:	1530
jouowing conumon notas.	1001
	1520

• $\exists \Phi \text{ such that } \Phi \models y, \text{ where } \Phi \text{ is a non-empty}$ 1532 subset of world knowledge \mathcal{K}_w and \models denotes 1533 semantic entailment. 1534

1535

1536

1537

1538

1539

1540

1541

Theorem 1 $\mathcal{F} \models \mathcal{T}, \mathcal{T} \not\models \mathcal{F}$, where \models denotes entailment.

Theorem 1 indicates that responses ensuring faithfulness are necessarily factual, but the converse may not always apply. The proof of the theorem is shown in the Appendix I.

I The proof of Theorem 1

Theorem. $\mathcal{F} \models \mathcal{T}, \mathcal{T} \not\models \mathcal{F}.$ 1542

	Please act as an impartial judge and evaluate the quality of the response provided
	by an AI assistant to the user utterance displayed below. You will be provided with
	the related knowledge, a dialogue history, and a generated response.
LLM-as-a-Judge	Your evaluation should be mainly based on <evaluation criteria=""></evaluation>
	Begin your evaluation by providing a short explanation. Be as objective as possible.
	After providing your explanation, you must rate the response on a scale of 1 to 5.
	The rating is:

Table 8: The prompts for LLM-based evaluations.

Context	<human> Madonna is one of my favorite entertainers, do you know much about her? <assistant> Sure! She was born in 1958 on August 16, and is an artist and an actress. <human> I didn't know she was that old, it doesn't matter though she's such a great artist anyway! <assistant> I see, well she is also cited by other artists as an influence. <human> I certainly think she's an influence as well, do you know any of Madonna's most famous songs?</human></assistant></human></assistant></human>
Knowledge	She followed it with a series of successful albums, including the global bestsellers, "Like a Virgin" (1984) and "True Blue" (1986), as well as the Grammy Award winners, "Ray of Light" (1998) and "Confessions on a Dance Floor" (2005).
Greedy	Yes! Some of her most famous songs are Material Girl, American Pie, Like a Virgin, and Ray of Light.
Nucleus	Some of her most famous songs include Like a Virgin, Material Girl, American Pie, and Hung Up.
DoLa	Of course! Some of her most famous songs include "Like a Virgin," "Material Girl," "Vogue," "Into the Groove," "Like a Prayer," "Express Yourself," "Cherish," "Frozen," "Hung Up," "4 Minutes"(More generated songs omitted here)
CAD	She has a series of successful albums, including the global bestsellers, "Like a Virgin" (1984) and "True Blue" (1986), as well as the Grammy Award winners, "Ray of Light" (1998) and "Confessions on a Dance Floor" (2005).
CoDe	Of course! She's had many! Like a Virgin, True Blue, Ray of Ligh, Confessions on a Dance Floor, are some of her most famous albums. Which song do you like the most?

Table 9: An example for case study. The song "American Pie" was covered by Madonna. The baseline confused song titles with album names and generated many songs that cannot be verified by external knowledge.

Context	<human> Are you a basketball fan?</human>
	<assistant> Yes, I am a fan of the five player sport. Are you?</assistant>
	<human> Not as much as I used to be. I watched The Bulls in the 90s when they were</human>
	the Dream Team.
Knowledge	Jordan played 15 seasons in the NBA for the Chicago Bulls and Washington Wizards.
Beam	Interesting! Michael Jordan played 15 seasons in the NBA for the Chicago Bulls and
	Washington Wizards.
Nucleus	Wow, impressive! Even though Michael Jordan played for several teams, he is mostly
	known for his time with the Chicago Bulls on the Dream Team.
CAD	I see. Jordan played 15 seasons in the NBA for the Chicago Bulls and
	Washington Wizards.
CoDe	Wow, that's amazing! Michael Jordan played for Chicago Bulls and Washington Wizards
	for 15 seasons. Do you like the Bulls because of Michael Jordan?

Table 10: An example for case study. Jordan only played for two teams, so it is inaccurate for Nucleus to say that he played for several teams. The "Dream Team" generally refers to the U.S. Olympic men's basketball team and should be distinguished from NBA games.

Instructions

You are given 3 comparison tasks. For each comparison task, you are given a dialogue context, a piece of external knowledge and two different knowledge-grounded responses.

For each dialogue, your task is to identify which of the **responses** is better based on naturalness, engagingness and faithfulness. If you aren't sure, select that they are **similar**.

Read the two knowledge-grounded responses of the dialogue context, and answer the 3 following questions about naturalness, engagingness and faithfulness:

Dialogue Context:

[Human]: Have you ever heard of Insane Clown Posse?

[Assistant]: Of course, its an American hip hop duo composed of Joseph Bruce and Shaggy 2 Dope.

[Human]: How familiar are you with the horrorcore genre they perform?

Related Knowledge:

Founded in Detroit in 1989, Insane Clown Posse performs a style of hardcore hip hop known as horrorcore and is known for its elaborate live performances.

Response A:

Insane Clown Posse is known for their horror-themed lyrics and elaborate live performances, making them pioneers of the horrorcore genre.

Response B:

I'm quite familiar with their horrorcore style. Their live performances are indeed quite elaborate and unique.

Naturalness (which response sounds more natural and fluent, matching human communication styles?):

- A is more natural.
- B is more natural.
- A and B are similar.

Creativenss (which response delivers knowledge interestingly and encourages you to keep interacting?):

A is more creative.

- B is more creative.
- A and B are similar.

Faithfulness (which response is consistent with the knowledge?):

- A is more faithful to the knowledge.
- B is more faithful to the knowledge.
- A and B are similar.

Figure 9: Human evaluation instructions and interface.

Context	<human> Hi, I like to volunteer at a nursing home, helping the elderly!</human>
Knowledge	Nursing homes will provide short-term rehabilitative stays following a surgery, illness
	or injury which may require physical therapy, occupational therapy or speech-language
	therapy.
Greedy	That's great to hear that you volunteer at a nursing home. Many nursing homes also
	offer long-term care for residents who need ongoing assistance with daily activities.
Nucleus	Great to hear that you're interested in volunteering at a nursing home!
CAD	Volunteering at a nursing home is a wonderful way to support the elderly. Many
	nursing homes offer short-term rehabilitative stays following surgeries, illnesses, or injuries,
	which may require physical therapy, occupational therapy or speech-language therapy.
CoDe	That's great! You can assist those who need physical, occupational, or speech therapy
	during their recovery. What a great job!

Table 11: An example for case study.

Proof. $\mathcal{F} \models \mathcal{T}$: For all y that satisfy $\mathcal{F}(K, y)$, there exists Γ such that $\Gamma \models y$ and $\Gamma \subseteq K$. Since $K \subsetneq K_w$ (external knowledge is a proper subset of world knowledge), it follows that $\Gamma \subseteq K_w$. Let $\Phi = \Gamma$, then $\Phi \models y$ and $\Phi \subseteq K$. Hence, $\mathcal{T}(y)$ holds, and the conclusion is proved.

 $\mathcal{T} \not\models \mathcal{F}$: We prove it by contradiction. Suppose that $\mathcal{T} \models \mathcal{F}$, then for all y that satisfy $\mathcal{T}(y)$, there exists Φ such that $\Phi \models y$ and $\Phi \subseteq K_w$. Let $\Phi \subseteq K_w/K$. Since $\mathcal{T} \models \mathcal{F}$, then $\Phi \subseteq K$. However, $\Phi \subseteq K_w/K$, it implies that $\Phi = \emptyset$, but $\Phi \neq \emptyset$, leading to a contradiction, thus the conclusion is not valid.