Effective Large Language Model Adaptation for Improved Grounding

Anonymous ACL submission

Abstract

001

011

012

015

022

034

042

Large language models (LLMs) have achieved remarkable advancements in natural language understanding and generation. However, one major issue towards their widespread deployment in the real world is that they can generate "hallucinated" answers that are not fac-Towards this end, this paper focuses tual. on improving LLMs by grounding their responses in retrieved passages and by providing citations. We propose a new framework, AGREE, Adaptation for GRounding EnhancEment, that improves the grounding from a holistic perspective. We start with the design of a test-time adaptation capability that can retrieve passages to support the claims that have not been grounded, which iteratively improves the responses of LLMs. To effectively enable this capability, we propose tuning LLMs to self-ground the claims in their responses and provide accurate citations. This tuning on top of the pre-trained LLMs requires well-grounded responses (with citations) for paired queries, for which we introduce a method that can automatically construct such data from unlabeled queries. Across five datasets and two LLMs, results show that our tuning-based AGREE framework generates superior grounded responses with more accurate citations compared to promptingbased approaches and post-hoc citing-based approaches.

1 Introduction

Recent advancements in large language models (LLMs) have yielded demonstrably groundbreaking capabilities in natural language processing (NLP) (Brown et al., 2020; Chowdhery et al., 2022). Their ability to understand, generate, and manipulate text at unprecedented scales and depths has established them as a transformative force within the burgeoning field of artificial intelligence, poised to significantly impact our increasingly datadriven world. Despite their widely spread adoption, one prominent issue of LLMs is that in certain scenarios they hallucinate: they generate plausiblesounding but nonfactual information (Maynez et al., 2020; Ji et al., 2023; Menick et al., 2022), limiting their the applicability in real-world settings. To mitigate hallucinations, solutions generally rely on grounding the claims in LLM-generated responses to supported passages by providing an attribution report (Rashkin et al., 2023; Bohnet et al., 2022; Gao et al., 2023a) or adding citations to the claims (Liu et al., 2023; Gao et al., 2023b; Huang and Chang, 2023). 043

045

047

049

051

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

There has been a growing amount of interest in making LLM-generated responses more trustworthy via grounding. One line of work follows the retrieval-augmented generation (Chen et al., 2017; Guu et al., 2020; Lewis et al., 2020) framework, in which LLMs are presented with retrieved passages or passages, and are instructed to include grounded responses in their answers via instruction tuning (Kamalloo et al., 2023) or in-context learning (Gao et al., 2023b; Kamalloo et al., 2023). As LLMs are required to perform this challenging task from just instructions or few-shot demonstrations, such directions often lead to mediocre grounding quality (Gao et al., 2023b). Another line of work is on post-hoc citing (Gao et al., 2023a; Chen et al., 2023), which links support passages to the claims in responses using an additional attribution evaluation model. This paradigm heavily relies on LLMs' parametric knowledge and might not scale to less known knowledge (Sun et al., 2023).

We propose a new framework, AGREE, Adaptation of LLMs for GRounding EnhancEment. As shown in Fig. 1, our framework adopts a learning-based approach to finetune LLMs, as opposed to relying on an external NLI model post-hoc. At the training phase, AGREE collects well-grounded responses for unlabelled queries automatically from a base LLM with the help of an attribution evaluation model (an



Figure 1: Our framework, AGREE, combines tuning and test time adaptation for better grounding.

NLI model). Next, the collected data is used for supervising LLMs to generate grounded responses based on the retrieved passages as well as include citations in their responses. At the testing phase, we propose an iterative inference strategy that allows LLMs to seek for additional information based on the self-grounding evaluation so as to refine its response. The tuning and test time adaptation together enable LLMs to effectively and efficiently ground their responses in the corpus.

087

101

102

103

104

105

106

107

109

110

111

We apply our AGREE framework to adapt an API-based LLM, text-bison, and an open LLM, 11ama-2-13b, with training data collected using unlabelled queries from three datasets. We conduct evaluation on both in-domain datasets and outof-distribution datasets and compare our AGREE framework against competitive in-context learning and post-hoc citing baselines. The experimental results highlight that AGREE framework successfully improves grounding (citation recall) and citation precision compared to the baselines by a substantial margin (generally more than 20%). We find LLMs can learn to add accurate citations to their responses with our carefully designed tuning mechanisms. Furthermore, the improvements in grounding quality achieved by tuning using certain datasets can generalize well across domains.

2 Related Work

112Hallucination is a prevalent issue for generative113language models on many tasks (Maynez et al.,1142020; Raunak et al., 2021; Dziri et al., 2021; Ji115et al., 2023; Tang et al., 2023; Huang and Chang,1162023), which leads to the development of system-117atic evaluation of the grounding in generated re-118sponses (Bohnet et al., 2022; Rashkin et al., 2023;

Min et al., 2023; Yue et al., 2023). Among these, our work focuses on providing citations to attributable information source (Liu et al., 2023; Gao et al., 2023b) for responses generated by LLMs. Unlike existing work that largely relies on zeroshot prompting or few-shot prompting (Kamalloo et al., 2023; Gao et al., 2023b), we use a learningbased approach that tunes LLMs to generate bettergrounded responses supported with citations. Furthermore, our approach teaches LLMs to cite the passages themselves as opposed to using an additional attribution evaluation model (Gao et al., 2023a; Chen et al., 2023). 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

Our framework improves LLMs for better grounding, as a form of a retrieval augmented generation approach. This differentiates work from recent work that improves factuality of LLMs without using retrieved passages by inference-time intervention (Li et al., 2023; Chuang et al., 2023), cross-exam (Cohen et al., 2023; Du et al., 2023), or self-verify (Dhuliawala et al., 2023), which cannot provide references. While past work have explored using retrieval to improve LLM generation quality (Chen et al., 2017; Lewis et al., 2020; Guu et al., 2020; Izacard and Grave, 2020; Shi et al., 2023) or factuality (Shuster et al., 2021; Jiang et al., 2023; Liangming Pan, 2023), our approach further uses self-generated attribution evaluation to provide citations and guide retrieval.

3 Problem & Background

We focus on enabling LLMs to provide grounded149responses – given an input text query Q and a corpus $\mathcal{D} = \{d_i\}$ consisting of text passages, we aim150to generate A to the query that is factually grounded152in the corpus \mathcal{D} . Our framework tunes a pre-trained153



Figure 2: Illustration of the proposed AGREE framework. A base LLM is tuned to generate grounded responses that include citations. At test time, self-generated grounding evaluation is used to iteratively improve the responses.

154 LLM $\mathcal{M}^{\mathcal{B}}$ to an adapted LLM $\mathcal{M}^{\mathcal{A}}$, primarily for 155 improving grounding with respect to the corpus 156 \mathcal{D} , which we evaluate using the grounding score 157 denoted as \mathcal{G} .

158

159

160

161

162

163

164

165

167

168

169

170

171

Grounding evaluation Let s_1, \ldots, s_n be the statements in the answer $A = s_1, \ldots, s_n$. In line with prior work (Rashkin et al., 2023; Gao et al., 2023a,b), we evaluate the grounding of an answer by assessing whether the statements in A can be attributed to the corpus \mathcal{D} . That is, we require the answers to be associated with citations $\mathcal{C} = \{E_i, \ldots, E_n\}$; each statement s_i to be linked a set of evidence passages $E_i \subset \mathcal{D}$. Then, the grounding quality of A can be quantified by:

$$\mathcal{G}(A, \mathcal{C}) = \frac{1}{n} \sum_{i} \phi(\operatorname{concat}(E_i), s_i)$$

where ϕ is an attribution evaluation model that assesses whether the concatenated passage concat(E_i) supports s_i .

4 AGREE Framework

The proposed AGREE framework takes a holistic 172 perspective for grounding, introducing a test-time 173 adaptation (TTA) mechanism to generate grounded 174 outputs that are supported from the corpus, and 175 proposing a model tuning approach to better align the pre-trained LLM with the desired TTA capabil-178 ity. In the following section, we first introduce the inference procedure of our method based on the 179 proposed iterative TTA approach, and then explain in detail how we tune the base LLM in order to enable such iterative TTA capability. 182

4.1 Test-time adaptation

At a high level, our framework is a form of retrieval augmented generation framework to output grounded responses. The inference procedure is overviewed in Algorithm (1). At the core of our approach lies an LLM that is able to answer a query based on a set of given passages retrieved from the corpus, and, more importantly, self-ground its response to add citations to the passages as well as to find unsupported statements needing further investigation. Using these capabilities, the LLM can guide the process of iteratively, constructing a set of relevant passages from the large corpus \mathcal{D} to refine its response to the query. 183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

As shown in Algorithm (1), given a query Q and the corpus \mathcal{D} , we first retrieve based on the query to obtain an initial set of working passages. Next, we employ the following procedure iteratively until we consume all the budget B of invoking LLM calls. At each iteration, the LLM generates a response to the query based on the working passages, adds citations to its response, and finds out any unsupported statements that do not have citations. Then, we add the cited passages to the list of relevant passages. Lastly, at each iteration, we update the working passages - if there are unsupported statements, we include additional information retrieved based on the unsupported statements (ln 11), otherwise, we include more passages that are retrieved based on the query to acquire more complete information regarding the query (ln 13). Note that at each iteration, we let the LLM to re-generate a response based on the current working passages instead of editing from previous one, which we observed lead to better fluency.

241

242

243

245

218

219

Algorithm 1 Iterative TTA

1:	procedure ITERATIVEINFERENCE($Q, \mathcal{D}, \mathcal{M}^{\mathcal{A}}, k, B$)
	input: A query Q, text corpus \mathcal{D} , the LLM for generating grounded response $\mathcal{M}^{\mathcal{A}}$, the number of passages k that $\mathcal{M}^{\mathcal{A}}$ can take as input, the budget for LLM
	calls B
2: 3: 4: 5: 6:	$relevant_psgs = []$ $ ightarrow retrieve passages using the query working_psgs := \text{RETRIEVE}(Q, D)[: k]while iter = 1 : B doightarrow$ Use the LLM to generate an answer A for the query Q based on the working psgs D. Additionally obtain the cited passages and unsupported the
	sentences.
7:	$A, cited_psqs, unsup_sents := \mathcal{M}^{\mathcal{A}}(Q, workinq_psqs)$
8:	▷ add cited passages to the list of relevant passages and de-duplicate the list
9:	$relevant_psgs := DEDUPLICATE(relevant_psgs + cited_psgs)$
10	if unsup_sents is not None then
11	: $working_psgs := DEDUPLICATE(relevant_psgs + Retrieve(unsup_sents, D))[:k]$
12	else
13	: $working_psgs := \text{DEDUPLICATE}(relevant_psgs + Retrieve(Q, D))[: k]$
14	: return A, cited_psgs

The design of our proposed TTA enables efficient and flexible inference. We rely on the LLM to generate citations itself, which has the advantage of reduced overhead of invoking another attribution evaluation model in a post-hoc way. Also, as we iteratively refine the answer, such a process can be streamed and flexibly controlled by setting a budget in deployment.

4.2 Tuning LLMs

Recall that our TTA requires the LLM to be able to self-ground its response, which we achieve by tuning the base LLM using data automatically created with the help of the attribution evaluation model.

Generating self-grounded responses We adapt the base LLM to generate grounded responses with citations. Our method is able to grant LLMs such an ability using only a collection of *unlabeled* queries Q and an attribution evaluation model ϕ . As we are using unlabeled queries, we formulate the adaptation task as tuning LLMs to achieve better grounding without heavily deviating from the original generations (this idea of preservation has also been adopted in recent work (Gao et al., 2023a)). Conceptually, we adapt $\mathcal{M}^{\mathcal{B}}$ to $\mathcal{M}^{\mathcal{A}}$ so that the answers generated by the adapted LLM $\mathcal{M}^{\mathcal{A}}$ should satisfy the grounding constraints (with grounding score > $\tau_{\mathcal{G}}$) while maximizing the scores with respect to the base LLM $\mathcal{M}^{\mathcal{B}}$:

 $\max \mathbb{E}_{A \sim \mathcal{M}^{\mathcal{A}}(\cdot | Q, \mathcal{D})} \mathcal{M}^{\mathcal{B}}(A \mid Q, \mathcal{D}) \mathbb{1}\{\mathcal{G}(A, \mathcal{C}) \geq \tau_{\mathcal{G}}\}$

247This leads to a data-centric approach for optimizing248 $\mathcal{M}^{\mathcal{A}}$. Since grounding score can vary with differ-249ent query characteristics (e.g., responses for more250open-ended questions are generally associated with251lower grounding scores compared to factoid ques-252tions), instead of discarding all data points with

grounding scores below a hard threshold, we instead encourage LLMs to generate a maximallygrounded response given a question. 253

254

255

256

257

258

259

260

261

262

263

264

265

267

268

269

270

271

272

273

274

275

276

277

278

279

281

282

283

284

287

Given the query, we first sample responses $\{A\}$ from the base LLM $\mathcal{M}^{\mathcal{B}}(\cdot \mid Q, \mathcal{D})$ using instruction following (see Appendix A for details). For each $A = s_1, \ldots, s_n$ we create citations $C = \{E_i\}$ using the attribution evaluation model, ϕ , to link a sentence s_i to the maximally supported passage $e_i = \max_{d \in \mathcal{D}} \phi(d, s_i)$ if the passage e_i actually support s_i (i.e., $\phi(e_i, s_i) > \tau$). Otherwise, we do not add a citation to s_i , and s_i is an unsupported statement. That is: $E_i = \{e_i\} i f \phi(e_i, s_i) >$ $\tau else \{\}$. We use U to denote the set of unsupported statements. This allows us to evaluate the grounding of A as in Section 3. Now, we can choose the best response A^* from $\{A\}$ based on the grounding scores to form a grounded response, i.e., $A^* = \arg \max_A \mathcal{G}(A, C)$.

We use $\{Q, A^*, C\}$ to teach the base LLM to generate grounded responses with citations. As shown in Fig. 2 (left), in addition to citations, we also instruct the LLM to clearly state the unsupported statements U. We note the tuning of framework does not force all training responses to be perfectly grounded. Instead, we supervise the LLM itself to identify unsupported statements. This allows the LLM to generate more flexibly and guide the retrieval process with its knowledge.¹

Supervised fine-tuning We have introduced how we construct supervision to instruct the LLM to add citations C and state unsupported statements U in its response. To effectively tune the LLM, we verbalize the entire process in natural language. We denote the verbalized natural language descrip-

¹Please refer to Appendix A for more details on the tuning method.

Dataset	Туре	Corpus	#					
	Trai	in						
NQ	Factoid QA	Wiki	2500					
StrategyQA	Multi-htop QA	Wiki	1000					
Fever*	Fact Checking	Wiki	1000					
	In-Distribution Test							
NQ	Factoid QA	Wiki	700					
StrategyQA	Factoid QA	Wiki	460					
	Out-of-Distribution Test							
ASQA	Ambiguous QA	Wiki	948					
QAMPARI	Multi-answer QA	Wiki	1000					
Enterprise	Customer Support QA	Enterprise	580					

Table 1: Statistics used for adaptation and test datasets. In addition to in-domain test datasets, we also investigate the generalization to out-of-distribution datasets that exhibit different reasoning processes or different corpus types.

tion as VERB (A^*, C, U) (see Fig. 2 for a concrete example). The natural language formalization also allows us to conveniently tune the LLM with standard language modeling objectives:

$$\mathcal{M}^{\mathcal{A}} = \operatorname*{arg\,max}_{\mathcal{M}} \sum_{Q} \mathcal{M}(\operatorname{VERB}(A^*, \mathcal{C}, U) \mid Q, \mathcal{D})$$

Training data We use multiple datasets to construct the adaptation data used to tune the pre-trained LLM, including Natural Questions (NQ) (Kwiatkowski et al., 2019), FEVER (Thorne et al., 2018), and StrategyQA (Geva et al., 2021). We choose these as they contain diverse text, and the answers to the corresponding queries require different types of reasoning processes: NQ provides diverse queries naturally asked by real users; FEVER places a particular emphasis on fact verification; and StrategyQA requires multi-hop reasoning with implicit strategy. It is worthwhile to note that AGREE *only* uses queries, leaving out ground-truth answers, to improve LLMs.

5 Experiments

5.1 Setup

Evaluation datasets We conduct comprehensive evaluation on 5 datasets. In addition to the two indomain test sets, NQ and StrategyQA (we leave 310 out the non-QA dataset, FEVER), we further test 311 the generalization of adapted LLMs on 3 out-of-312 domain datasets, including ASQA (Stelmakh et al., 314 2022), QAMPARI (Amouyal et al., 2022), and an Enterprise dataset. In particular, ASQA and QAM-315 PARI contain questions of ambiguous answers and multiple answers. The Enterprise dataset is a proprietary dataset which requires provided answers 318

that are grounded in customer service passages. Such an evaluation suite allows assessing the generalization capability of the adapted LLMs for OOD question types (ASQA and QAMPARI) as well as to an entirely different corpus (Enterprise). 319

320

321

322

323

324

325

326

328

329

330

331

332

333

334

335

336

337

338

341

342

343

345

346

347

348

349

350

351

353

354

355

356

357

358

359

360

361

362

363

364

365

Models We demonstrate AGREE framework with two LLMs, text-bison and LlaMA-2-13B (Touvron et al., 2023). We use GTR-large (Ni et al., 2021) as our retriever, and use TRUE (Honovich et al., 2022) as the attribution evaluation model.

Baselines We evaluate the effectiveness AGREE in two settings, invoking LLMs once, without TTA; and invoking LLMs multiple times, with the proposed TTA. We compare with three baselines from recent work, including one prompting-based approach and two post-hoc citing approaches, described below.

Few-shot In-Context Learning (ICLCITE): Following Gao et al. (2023b), we prompt LLMs with few-shot examples (Gao et al., 2023b), each consisting of a query, a set of retrieved passages, and an answer with inline citations. The LLMs can therefore learn from the in-context examples and generated citations in the responses. It is worth-while to note that **ICLCITE do have access to retrieved passages**.

Post-hoc search (POSTSEARCH): Following Gao et al. (2023b), given a query, we first instruct LLMs to answer the query *without* passages, and then add citations in a post-hoc way via searching. We link each claim in the response to the most relevant passage retrieved from a set of query-related passages. This baseline only uses retriever but not the attribution model, ϕ .

Post-hoc Attribution (POSTATTR): Following Gao et al. (2023a), instead of citing the most relevant passage, for each claim, we retrieve a set of k passages from the corpus, and then use the attribution evaluation model to link to the passage that maximally supports the claim. We note both baselines in the post-hoc citing paradigm only rely on LLMs' parametric knowledge.²

Metrics We mainly focus on improving the grounding quality of generated responses, reflected by the quality of citations. Following past work (Gao et al., 2023b), we report the **citation recall** (rec) and **citation precision** (pre) for all the

294

295

303

305

²Please refer to Appendix B for more details on experimental setup.

	NQ		StrategyQA			ASQA			QAMPARI			Enterprise		
	em-rec	rec	pre	acc	rec	pre	em-rec	rec	pre	rec-5	rec	pre	rec	pre
						Base	model: text-bison-00			1				
ICLCITE	47.6	52.1	56.3	74.5	13.6	27.8	39.5	47.3	49.8	20.3	22.7	24.5	30.2	40.5
POSTSEARCH	45.1	29.7	28.7	75.5	20.1	20.1	38.4	19.2	19.2	22.5	16.2	16.2	15.9	15.9
PostAttr	45.1	31.5	31.5	75.5	18.4	18.4	35.1	38.0	38.0	22.5	18.5	18.5	20.1	20.1
AGREE _{W/O TTA}	50.0	67.9	73.1	74.1	33.4	50.5	39.5	65.9	70.5	20.1	60.1	64.5	55.8	67.1
AGREE _{w/ TTA}	53.1	70.1	75.0	74.9	39.2	57.9	40.9	73.2	77.0	20.9	62.9	67.1	57.2	68.6
						Bas	e model: 1	lama-	-2-13b					
ICLCITE	45.8	42.8	41.6	65.5	20.6	33.1	35.2	38.2	39.4	21.0	10.2	10.4	30.6	38.8
POSTSEARCH	35.9	17.5	17.5	64.3	8.7	8.7	25.0	23.6	23.6	12.0	27.5	27.5	13.4	13.4
PostAttr	35.9	26.0	26.0	64.3	12.5	12.5	25.0	33.6	33.6	12.0	28.9	28.9	18.7	18.7
AGREE _{w/o TTA}	47.9	50.5	56.6	65.0	25.5	35.0	35.7	50.2	55.3	17.1	40.4	43.6	50.6	53.8
AGREE _{W/ TTA}	51.0	62.0	66.0	64.6	30.2	37.2	39.4	64.0	66.8	17.9	51.4	53.4	50.4	55.4

Table 2: Answer accuracy and grounding (measured by citation quality) of AGREE and baselines across 5 datasets. Our approach achieves substantially better citation grounding (measured by citation recall) and citation precision compared to the baselines.

datasets we are evaluating on. We note that **cita**tion recall aggregates how well each sentence is supported by the citation to the corpus, which is essentially the grounding score \mathcal{G} . Therefore, we prioritize on the evaluation citation recall.

We also report the correctness of the generated outputs. For NQ, we report exact match recall (emrec; whether the short answers are substrings in the response). For StrategyQA, we report the accuracy (acc). For ASQA and QAMPARI, we use subsets from Gao et al. (2023b), and report the exact match recall (em-rec) for ASQA and recall-5 (rec-5, considering recall to be 100% if the prediction includes at least 5 correct answers) for QAMPARI. For the Enterprise dataset, we only report the citation quality as there are no ground truth answers for this dataset, and citation quality reflects whether the model can provide accurate information.

5.2 Main results

367

371

374

376

384

390

394

395

398

Tuning is effective for superior grounding Table 2 summarizes the results obtained using our AGREE framework and compares with the baselines. As suggested by the results, across 5 datasets, AGREE can generate responses that are better grounded in the text corpus and provide accurate citations to its response, substantially outperforming all the baselines. When tuned with high-quality data, LLMs can effectively learn to self-ground their response without needing an additional attribution model. By contrast, ICLCITE, which solely relies on in-context learning, cannot generate citations as accurately as a tuned LL, as suggested by the large gap on citation precision between ICLCITE and AGREE. We also observe similar findings as suggested by Gao et al. (2023b): POSTCITE often leads to poor citation quality – without being conditioned on passages, the response from POSTCITE often cannot be paired with passages that lead to high attributions score for the generated claims. 399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

The performance improvements can generalize Recall that we adapt the base LLM only using in-domain training sets (NQ, StrategyQA, and FEVER), and directly test the model on out-ofdistribution (OOD) test set (ASQA, QAMPARI, Enterprise). The results suggest that the improvements obtained from training on in-domain datasets can effectively generalize to OOD datasets that contain different question types or use different types of corpus. This is a fundamental advantage of the proposed approach – AGREE can generalize to a target domain in the zero-shot setting without needing any samples from the target domain, which is needed for ICLCITE.

TTA improves both grounding and answer correctness The comparison between AGREE without and with TTA highlights the effectiveness of our iterative TTA strategy. We observe improvements in terms of both better grounding and accuracy. For instance, TTA improves 11ama-2 answer correctness by 3.1 and 3.7 on NQ and ASQA, respectively. Such improvements can be attributed to the fact that our TTA allows the LLMs to actively collect relevant passages to construct better answers following the self-grounding guidance.

		NQ		StrategyQA		ASQA			QAMPARI			Enterprise		
	em-rec	rec	pre	acc	rec	pre	em-rec	rec	pre	rec-5	rec	pre	rec	pre
						Base	model: te	xt-bis	son-001					
ICLCITE	47.6	52.1	56.3	74.5	13.6	27.8	39.5	47.3	49.8	20.3	22.7	24.5	30.2	40.5
AGREEW/0 TTA	50.0	67.9	73.1	74.1	33.4	50.5	39.5	65.9	70.5	20.1	60.1	64.5	55.8	67.1
AGREE _{W/0 TTA}	49.4	62.3	69.1	74.1	33.0	45.5	38.4	56.0	64.5	19.1	43.7	49.5	40.5	59.2
					Base model: 11ama-2-13b									
ICLCITE	45.8	42.8	41.6	65.5	20.6	33.1	35.2	38.2	39.4	21.0	10.2	10.4	30.6	38.8
AGREEW/0 TTA	47.9	50.5	56.6	65.0	25.5	35.0	35.7	50.2	55.3	17.1	40.4	43.6	50.6	52.8
AGREE _{W/0 TTA}	48.1	47.4	53.6	62.1	25.0	30.2	35.0	44.0	51.2	15.7	33.1	38.0	44.7	49.2
AGREE ^{Distill} W/O TTA	47.9	59.1	65.1	64.4	30.5	41.1	35.2	58.5	65.2	17.9	52.5	52.7	48.1	55.9

Table 3: Analysis on the impact of training data. Training with multiple datasets (AGREE^{Multi-dataset}) leads to better grounding (citation recall) and better citation precision across datasets, compared to training using the NQ dataset (AGREE^{NQ-only}). The citation quality of a less capable model llama-2-13b can also benefit from tuning using outputs from a more capable model (text-bison-001).

	# Tok: LLM	# Tok: NLI (T5-11B)
ICLCITE PostAttr	2800 360	3520
AGREE _{w/o TTO} AGREE _{w/ TTO}	1210 4840	

Table 4: The average computation cost (for one query) of different methods measured by number of tokens processed by the LLM and the NLI model (a T5-11B model). AGREE_{w/o TTO} is able to achieve better citation quality compared to ICLCITE, despite consuming less than half of the tokens needed for ICLCITE.

Discussion on answer correctness In general, 431 AGREE_{w/TTO} can achieve better correctness com-432 pared to ICLCITE. AGREE_{w/o TTO} achieves similar 433 answer correctness with ICLCITE, as both meth-434 ods are conditioned on the same set of passages. As 435 a result, the quality of passages heavily intervenes 436 on the correctness of the answers. Unlike AGREE 437 and ICLCITE, POSTATTR purely relies on the para-438 439 metric knowledge of the LLMs to answer the query. As a result, POSTATTR generally achieves infe-440 rior answer correctness compared to AGREE and 441 ICLCITE on these two LLMs, especially on the 442 less capable LLM, 11ama-2-13b, that has less ac-443 curate knowledge compared to bison. Moreover, 444 on the Enterprise dataset which contains very spe-445 cific information, POSTATTR utterly fails to recall 446 attributable information from LLMs' parametric 447 knowledge. 448

449Discussion on LLMsOur approach successfully450adapts both text-bison-001 and llama-2-13b. llama451is generally less capable compared to bison, un-452derperforming bison in terms of answer correct-453ness and citation quality. Still, AGREE also consis-454tently outperforms the baseline, generating more

grounded answers as well as providing more precise citations. This shows our tuning-based adaptation is model-agnostic and is effective across LLMs of varying capabilities. 455

456

457

458

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

484

485

486

487

488

489

5.3 Analysis

Efficiency Our AGREE framework finetunes the base LLM to enable self-grounding without needing for additional in-context examples or attribution model. As a result, our framework is able to achieve strong citation performance without expensive inference cost. 2 Table 4 shows the comparison between the computation cost, measured by the number of tokens processed by the LLM and attribution model, needed for one query of our methods and that of the baselines. Compared to ICLCITE, AGREE_{w/o TTO} uses much fewer tokens due to not using additional in-context examples, but achieves significantly better citation quality (see Table 2). POSTATTR does not use retrieved passages in the prompts and hence requires less computation on the LLM compared to our framework, but it requires additional overhead of extensively invoking the NLI model (which is also of 11B parameters, a relatively large scale) to verify the each of claim based on each of the retrieved passages. The citation performance of POSTATTR also substantially lags ICLCITE and AGREE. AGREE_{w/TTO} requires more computation compared to AGREE_{w/o TTO}, but is able to achieve both better citation quality and improvements in answer correctness.

Impact of multi-dataset training Our AGREE framework multiple datasets spanning factoid QA, multi-hop reasoning, and fact-checking to construct data for adapting the base model. We expect such a combination can grant the adapted model better

generalization to different types of questions and 490 different text distribution. We conduct an analysis 491 to investigate the benefits of using multiple datasets 492 for tuning. Table 3 shows the performance of our 493 approach trained using multi-datasets and a counterpart that is trained only on NQ data (AGREE^{NQ-only}). 495 The results suggest that training using NO leads 496 to inferior citation quality compared to training 497 on the combination of three datasets across all the 498 datasets. The performance gap is especially signifi-499 cant on datasets other than NQ. Moreover, training on NQ only also leads to inferior answer correct-501 ness across all the datasets, which also suggests the 502 benefits of using multiple datasets. Nevertheless, training only on NQ can still improve the perfor-504 mance compared to solely relying on in-context learning (ICLCITE).

Distilling from bison to llama-2-13b Our work mainly focuses on improving the base LLM's grounding capability in a self-improving fashion. That is, we use the samples generated by the base 510 LLM itself to adapt the base LLM, as opposed 511 to distilling from proprietary models which may impose constraints on the deployment of adapted 513 model. Nevertheless, we conduct an analysis to 514 investigate the effectiveness of distilling the data 515 516 produced by a more capable model to enhance the grounding of a less capable. In our case, we 517 use the data generated by text-bison-001 to tune 518 llama-2-13b. As shown in the last row of Table 3, 519 AGREE^{Distill} achieves better citation quality com-520 pared to AGREE on llama-2-13b, as it is trained on better grounded responses produced by a more ca-522 pable model. However, 11ama-2 finetuned on data generated by bison still can't level the performance 524 on bison as constrained by the gap between the capabilities of these two LLMs.

Qualitative analysis We qualitatively analyze 527 the advantages of proposed AGREE framework compared to ICLCITE, the strongest among the 529 baselines. We observe that on both text-bison-001 530 and 11ama-2-13b, ICLCITE achieves inferior citation quality due to failure in following the citation format (e.g., adding citations after the periods, violating the instructions), linking a statement to a relevant but un-attributable passage (as indi-536 cated by poor citation precision), and introducing more auxiliary information not mentioned in the retrieved passages (as indicated by citation recall). Our AGREE framework mitigates these issues by tuning on well-grounded responses certified by the 540



Figure 3: Output examples of the proposed AGREE framework with text-bison-001 as the base model. TTA is able to improve the response by retrieving more relevant information to precisely support a statement (see top) or finding more passages to generate a more complete response (see bottom).

attribution evaluation model. We also provide example outputs in Fig. 3 comparing the outputs of AGREE with and without proposed TTA and observe that TTA can help find more supporting passages by active retrieving using unsupported statements (top) or iteratively find more passages to construct a more complete response (bottom).

6 Conclusion

We have introduced a novel framework, AGREE, that adapts LLM for improved grounding. Our framework tunes a pre-trained LLM to self-ground its response in retrieved passages using automatically collected data. The integrated capability for grounding their responses further enables the LLM to improve the responses at test time. Our evaluations across five datasets demonstrate the benefits of the proposed learning-based approach compared to approaches that solely rely on prompting or the parametric knowledge of LLMs.

558

559

541

7 Limitations

560

561

562

564

568

571

573

575

576

580

584

585

586

590

591

596

597

606

607

610

Our approach employs an automated data creation method that relies on an attribution evaluation model to create citations instead of hiring humans to annotate citations. Thus, the citation quality is dependent on the performance of the attribution evaluation model. As suggested in Gao et al. (2023b); Honovich et al. (2022), the model we use favors "fully support" and cannot effectively detect "partially support". The adapted LLMs may favor adding "fully support" citations as a result. One solution is to curate a set of human-annotated citations for "partially support", which we defer to future work.

Our evaluation follows prior work (Rashkin et al., 2023; Gao et al., 2023a) and uses the attribution evaluation model to evaluate grounding and citation quality. Therefore, our work can encounter the same issue as in past work: the citation grounding and citation quality evaluation is limited by the capability of the NLI model.

Our approach uses created grounded responses to LLMs via supervised finetuning, as we observed this straightforward tuning technique leads to empirical strong performance. It is also possible to treat grounding as a preference and RLHF (Ouyang et al., 2022) to tune LLMs, which we leave as future work.

This work focuses only considers open domain question answering datasets focusing on information seeking and written in the English language. It is unsure how well the framework can handle other language.

Lastly, this work studies the approach of adding citations to LLM-generated response and which carries a shared risk with related research: a seemingly plausible but incorrect citation could potentially make an unsupported statement more convincing to users.

References

- Samuel Joseph Amouyal, Tomer Wolfson, Ohad Rubin, Ori Yoran, Jonathan Herzig, and Jonathan Berant. 2022. Qampari: An open-domain question answering benchmark for questions with many answers from multiple paragraphs.
- Bernd Bohnet, Vinh Q Tran, Pat Verga, Roee Aharoni, Daniel Andor, Livio Baldini Soares, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, et al. 2022. Attributed question answering: Evaluation and modeling for attributed large language models. *arXiv preprint arXiv:2212.08037*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. 611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

- Anthony Chen, Panupong Pasupat, Sameer Singh, Hongrae Lee, and Kelvin Guu. 2023. Purr: Efficiently editing language model hallucinations by denoising language model corruptions. *arXiv preprint arXiv:2305.14908*.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to answer opendomain questions. In Association for Computational Linguistics (ACL).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Baindoor Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Oliveira Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways. ArXiv. abs/2204.02311.
- 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. *arXiv preprint arXiv:2309.03883*.
- Roi Cohen, May Hamri, Mor Geva, and Amir Globerson. 2023. Lm vs lm: Detecting factual errors via cross examination. *ArXiv*, abs/2305.13281.
- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. Chain-of-verification reduces hallucination in large language models. *ArXiv*, abs/2309.11495.

777

778

724

Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate. *ArXiv*, abs/2305.14325.

670

671

688

697

705

706

707

710

711

712

713

714

715

716

717

718

719

721

722

- Nouha Dziri, Andrea Madotto, Osmar Zaïane, and Avishek Joey Bose. 2021. Neural path hunter: Reducing hallucination in dialogue systems via path grounding. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2197–2214, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. 2023a. RARR: Researching and revising what language models say, using language models. In *Proceedings of the Annual Conference of the Association for Computational Linguistics (ACL).*
 - Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023b. Enabling large language models to generate text with citations. In *Proceedings of the Conference* on Empirical Methods in Natural Language Processing (EMNLP).
 - Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929– 3938. PMLR.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*.
- Or Honovich, Roee Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022. TRUE: Re-evaluating factual consistency evaluation. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3905–3920, Seattle, United States. Association for Computational Linguistics.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Jie Huang and Kevin Chen-Chuan Chang. 2023. Citation: A key to building responsible and accountable large language models. *ArXiv*, abs/2307.02185.

- Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering.
- 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 Comput. Surv.*, 55(12).
- Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Li-Yu Daisy Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. *ArXiv*, abs/2305.06983.
- Ehsan Kamalloo, Aref Jafari, Xinyu Crystina Zhang, Nandan Thakur, and Jimmy Lin. 2023. Hagrid: A human-llm collaborative dataset for generative information-seeking with attribution. *ArXiv*, abs/2307.16883.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Inferencetime intervention: Eliciting truthful answers from a language model.
- Xinyuan Lu Anh Tuan Luu William Yang Wang Min-Yen Kan Preslav Nakov Liangming Pan, Xiaobao Wu. 2023. Fact-checking complex claims with program-guided reasoning. In *Proceedings* of the 61th Annual Meeting of the Association for Computational Linguistics (ACL 2023), Toronto, Canada.
- Nelson F. Liu, Tianyi Zhang, and Percy Liang. 2023. Evaluating verifiability in generative search engines. ArXiv:2304.09848.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1906–1919, Online. Association for Computational Linguistics.
- Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick,

- 779 780
- 70 78
- 78
- 78
- 78 78
- 78
- 790
- 79

803

805

808

809

811

812

813

814

815

816

817

818

819

823

824

826

827

831

835

- Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, and Nathan McAleese. 2022. Teaching language models to support answers with verified quotes. *ArXiv*, abs/2203.11147.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023.
 Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei Yang. 2021. Large dual encoders are generalizable retrievers. *ArXiv*, abs/2112.07899.
 - Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. 2022. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155.
 - Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm, Lora Aroyo, Michael Collins, Dipanjan Das, Slav Petrov, Gaurav Singh Tomar, Iulia Turc, and David Reitter. 2023. Measuring attribution in natural language generation models. *Computational Linguistics*, pages 1–64.
 - Vikas Raunak, Arul Menezes, and Marcin Junczys-Dowmunt. 2021. The curious case of hallucinations in neural machine translation. In *Proceedings of the* 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1172–1183, Online. Association for Computational Linguistics.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023. Replug: Retrievalaugmented black-box language models. ArXiv, abs/2301.12652.
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation. In *Findings* of the Association for Computational Linguistics: *EMNLP 2021*, pages 3784–3803, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022. ASQA: Factoid questions meet long-form answers. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 8273–8288, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Kai Sun, Y. Xu, Hanwen Zha, Yue Liu, and Xinhsuai Dong. 2023. Head-to-tail: How knowledgeable are large language models (llm)? a.k.a. will llms replace knowledge graphs? *ArXiv*, abs/2308.10168.

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

869

- Liyan Tang, Tanya Goyal, Alexander R. Fabbri, Philippe Laban, Jiacheng Xu, Semih Yavuz, Wojciech Kryściński, Justin F. Rousseau, and Greg Durrett. 2023. Understanding factual errors in summarization: Errors, summarizers, datasets, error detectors. In *Proceedings of ACL*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- Xiang Yue, Boshi Wang, Kai Zhang, Ziru Chen, Yu Su, and Huan Sun. 2023. Automatic evaluation of attribution by large language models. *arXiv preprint arXiv:2305.06311*.

A Details of Tuning

871

879

881

Recall that we create the tuning data by first sampling responses from the base LLM and then using
the attribution evaluation model to create citations
and identify unsupported statements. We will detail
the process in the following of this section.

Corpus & retriever As mentioned before, our framework is an instantiation of retrievalaugmented framework. For the datasets using Wikipedia as the corpus (NQ, StrategyQA, ASQA, and Qampari), we use the 2018-12-20 Wikipedia snapshot as the corpus and set up the retriever using GTR-large (Ni et al., 2021).

Task: You will be given a question and some search results. Please answer the question in 3-5 sentences, and make sure you mention relevant details in the search results. You may use the same words as the search results when appropriate. Note that some of the search results may not be relevant, so you are not required to use all the search results, but only relevant ones.

<Question>

Search Results: [<Index>] <Title> <Text>

[...]

Answer:

Figure 4: Zero-shot prompt template for sampling initial responses from the base LLM.

Sampling initial responses We sample initial responses from the base LLM using instruction following in a *zero-shot fashion*. Given a query, we present the base LLM with query and retrieved passages appended after an instruction that requires the base LLM to answer the query based on the passages; see Fig. 4 for the template of the zeroshot prompt. We note that we opt to use a zeroshot prompt as opposed to a task-specific few-shot prompt since 1) this can avoid biasing the generation with the few-shot in-context examples, and 2) this matches the expected scenario for deploying the adapted LLM to handle new queries in a zero-shot fashion.

For text-bison-001, we sample 4 responses using a temperature of 0.5. For 11ama-2-13b, we sample 4 responses using nuclear sampling (Holtzman et al., 2019) with p=0.95.

Adding citations and identifying unsupporteddocumentsAfter obtaining the initial response

Input

Task: You will be given a question and some search results. You are required to perform the following steps.

First, please answer the question in 3-5 sentences, and make sure you mention relevant details in the search results. You may use the same words as the search results when appropriate. Note that some of the search results may not be relevant, so you are not required to use all the search results, but only relevant ones. If you use the provided search results in your answer, add [n]-style citations.

Next, review your response and find the unsupported sentences that do not have citations.

<Question>

Search Results: [<Index>] <Title> <Text>

[...]

Output

Answer: <Response with citations>

Sentences Not Supported by Citations: <Unsupported statements>

Figure 5: Verbalization template for creating the training data for adapting the base LLM.

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

{*A*} from the base LLM. We break each response *A* into sentences into s_1, \ldots, s_i . For each s_i , we find the maximally supported passage e_i (scored by $\phi(e_i, s_i)$) that the base LLM has seen during generating the initial responses. We link e_i to s_i if $\phi(e_i, s_i) > 0.7$ to encourage more precise citations. For a sentence s_i if there does *not* exist an e_i such that $\phi(e_i, s_i) > 0.5$ (the decision boundary for entailment), we add s_i to the unsupported statement set *U*.

Verbalizing We show the template for verbalizing the data used to tune the LLM in Fig. 5. As shown in the figure, we verbalize the citations in enclosed box brackets that are added at the end of sentences (before periods) like [n], and verbalize unsupported statements after the responses.

B Details of Experiments

For tuning, we use LORA tuning (Hu et al., 2022) in experiments on both text-bison-001 and 11ama-2-13b. For bison, we use API to perform tuning.³ and follow all the default hyper-parameters except for training steps. We set 10% data created

³https://cloud.google.com/vertex-ai/docs/generativeai/models/tune-text-models-supervised

as development data and choose to use a training 926 step of 1000 (chosen from 500, 1000, and 2000). 927 For 11ama-2, we use the huggingface transform-928 ers (Wolf et al., 2019) chat-version checkpoint.⁴ 929 We find the chat-version achieves better perfor-930 mance than the base checkpoint in our prelimi-931 nary investigation. We set lora r to be 32, and 932 only choose to use a learning rate of 1e-5 (chosen from 1e-4 and 1e-5) using the development set. We finetune 11ama-2 on two A100 (40GB) GPU for 4 935 epochs. 936

937

939

941

942

943

944

945

946

947

951

952

953

957

Our evaluation uses official code from ALCE (Gao et al., 2023b), we use the same data split and prompt template from ALCE. We use temperature 0.25 for evaluation on both bison and llama. We use one sample for evaluation since adapted LLMs tend to generate better-grounded response exhibiting less variation.

C Comparison to ICLCITE on More Capable LLMs

	А	SQA		QAMPARI					
	em-rec E	rec Base m	pre odel: 1	rec-5 lama-2	rec -13b	pre			
AGREE _{W/O TTA}	35.7	50.2	55.3	17.1	40.4	43.6			
$AGREE_{W/\ TTA}$	39.4	64.0	66.8	17.9	51.4	53.4			
	Base model: 11ama-2-70b								
ICLCITE	41.5	62.9	61.3	21.8	15.1	15.6			
	Base model: ChatGPT-0301								
ICLCITE	40.4	73.6	72.5	20.8	20.5	20.9			

Table 5: Comparing AGREE on llama-2-13B against ICLCITE on llama-2-70B and ChatGPT-0301. We directly quote results from ALCE.

Table 5 compares AGREE using 11ama-2-13B as the base model against ICLCITE on more capable models. We directly use the results from ALCE (Gao et al., 2023b). Our framework is able to substantially shorten the gap between a small llama-2 model and much more capable LLMs.

D License of Datasets

The licenses datasets used in our work include:

- NQ (Kwiatkowski et al., 2019) under Creative Commons Share-Alike 3.0 license.
- StrategyQA (Geva et al., 2021) under MIT License.

- Fever (Thorne et al., 2018) under Creative 958 Commons Share-Alike license. 959
 Ambiguous QA (Stelmakh et al., 2022) under Creative Commons Share-Alike 3.0 license. 961
- Qampari (Amouyal et al., 2022) under Creative Commons Zero v1.0 Universal license.

⁴https://huggingface.co/meta-llama/Llama-2-13b-chat-hf