
CREAM-RAG: Enhanced Retrieval Augmented Generation to Limit Hallucination through Consistency-based Self-RAG

Yuliah Louis, Vivek Sekhadia, James Vaisman, Kingston Huynh, Sri Yanamandra¹, Kevin Zhu², Ryan Lagasse²

¹University of Illinois Urbana–Champaign ²Algoverse AI

kevin@algoverseairsearch.org, ryan@algoverseairsearch.org

Abstract

1 Retrieval-Augmented Generation (RAG) systems enhance large language models
2 (LLMs) by grounding responses in external evidence, mitigating hallucinations,
3 and enabling access to up-to-date, domain-specific knowledge. However, existing
4 RAG frameworks often suffer from unstable self-supervised optimization signals
5 and inconsistent factual grounding. We introduce CREAM-RAG (Consistency-
6 Regularized Enhanced Augmented Model for RAG) (Wang et al. [2025c]), a unified
7 framework that integrates retrieval, Direct Preference Optimization (DPO)-based
8 self-reward reinforcement learning, and a consistency regularization objective to
9 stabilize reward dynamics during fine-tuning. By enforcing alignment between
10 multiple retrieved contexts and generated responses, CREAM-RAG improves
11 factual faithfulness and semantic coherence without external supervision. Empirical
12 evaluations on the LLaMA-2-7B model demonstrate that CREAM-RAG achieves a
13 35.04% average improvement over the base model across reasoning and factuality
14 benchmarks, highlighting its effectiveness in reducing hallucinations and enhancing
15 retrieval-grounded reasoning.

16 1 Introduction

17 Retrieval-Augmented Generation (RAG) enhances language models by grounding outputs in retrieved,
18 relevant documents, effectively addressing key limitations such as hallucinations, updated knowledge,
19 and restricted domain expertise (Lewis et al. [2021]; Kwiatkowski et al. [2019]). By integrating
20 external sources during the inference time, RAG provides accurate and up-to-date facts. Recent
21 advancements, including autonomous retrieval frameworks, further improve efficiency and scalability
22 by reducing dependency on larger context windows and lowering computational costs (Zhou and
23 Chen [2025]).

24 Despite these improvements, existing self-rewarding RAG approaches, which focus on refining
25 retrieval and reward design, overlook a critical weakness: the inherent instability of self-generated
26 reward signals during training. Such systems are prone to failure modes like reward hacking and
27 retrieval-blind collapse, with recent work identifying temporal inconsistency as a fundamental yet
28 unsolved issue (Niu et al. [2024]).

29 We propose stabilizing the self-reward process in RAG systems to reduce hallucinations and reward
30 misalignment. Rather than redesigning reward types, our method enhances their reliability over
31 time by applying consistency regularization to self-reward signals, drawing on recent advances from
32 Consistency Regularized Self-Rewarding Language Models (CREAM) (Wang et al. [2025c]).

33 We introduce CREAM-RAG, which combines a retrieval module with an actor-critic generator where
34 the critic assigns self-rewards. Crucially, we augment standard Direct Preference Optimization
35 (DPO) with a consistency loss that minimizes divergence in reward signals across training steps. By
36 incorporating a frozen reference model to penalize shifts in reward judgments, our approach ensures
37 more stable and trustworthy reinforcement signals throughout the generation process.

38 **2 Related Works**

39 Our method integrates advances in RAG, self-rewarding RL, and consistency regularization to reduce
40 hallucinations and improve reliability.

41 **2.1 Other Variants to RAG and Early RAG**

42 Early RAG models showcased how coupling large language models with external retrieval can
43 improve accuracy. For example, REALM introduced an end-to-end retriever-generator framework
44 that significantly boosted open-domain question and answering (QA) performance and interpretation
45 capabilities (Guu et al. [2020]). Building on this, newer variants such as PIKE-RAG adapt retrieval
46 and rationale generation to better serve specialized applications (Wang et al. [2025a]). At the same
47 time, systematic surveys of RAG methods emphasize both the versatility and persistent challenges of
48 balancing retrieval efficiency and robustness (Oche et al. [2025]).

49 More recent work has focused on addressing two major challenges. Self-RAG and other self-
50 rewarding methods train models to adaptively decide when to retrieve, how to generate, and even how
51 to critique their own outputs, though they remain vulnerable to noisy or limited evidence (Asai et al.
52 [2023]). On the other hand, consistency-regularization approaches like CORD encourage models to
53 produce stable outputs under retrieval uneasiness while accounting for passage ranking, but often
54 strike the balance poorly, either ignoring or over-emphasizing rank information (Lee et al. [2024]).

55 Our work integrates core principles from RAG, self-rewarding reinforcement learning, and con-
56 sistency regularization. Existing RAG frameworks have enhanced factual accuracy by leveraging
57 external knowledge, yet they frequently exhibit unstable reward dynamic during reinforcement learn-
58 ing. For instance, while methods like Self-RAG (Asai et al. [2023]) incorporate self-assessment
59 mechanisms, they do not enforce consistency across training steps, leaving them susceptible to reward
60 hacking and preference drift. Conversely, consistency-focused approaches such as CORD (Lee et al.
61 [2024]) promote generation stability but fall short in merging retrieval awareness with self-rewarding
62 objectives.

63 CREAM-RAG bridges these lines of research by embedding consistency-regularized self-reward
64 mechanisms within the RAG process. Central to our approach is the use of a frozen reference model to
65 stabilize self-reinforcing RAG systems. By integrating this consistency mechanism into an actor-critic
66 RAG architecture and evaluating it across diverse tasks - including long-form QA, comprehension,
67 and hallucination mitigation - we show that maintaining internal consistency is essential for reducing
68 hallucinations and improving factual reliability. Our framework thus unifies previously separate
69 advances in retrieval quality, self-supervision, and reward stability.

70 **2.2 RAG Systems**

71 RAG enhances factuality by retrieving external documents during inference time (Lewis et al. [2021]).
72 Early systems used loosely coupled retrieval and generation, often yielding ungrounded outputs.
73 Recent training approaches optimize the retriever and generator, although noisy retrieval can still
74 cause hallucinations (Zhou and Chen [2025]). Self-rewarding reinforcement learning methods aim
75 to improve outputs, but suffer from unstable rewards and preference drift (Wang et al. [2025b]).
76 Consistency regularization helps stabilize training by penalizing shifts in reward preferences (Zhang
77 et al. [2024b]).

78 **2.3 Long-form QA and Challenges**

79 Long-form QA introduces additional challenges for RAG, requiring synthesis across lengthy and con-
80 flicting documents. Early solutions either overloaded context or oversimplified outputs (Kwiatkowski
81 et al. [2019]). Recent methods use aspect-based summarization (Hayashi et al. [2020]) and modular

82 document processing (Izacard and Grave [2021]), with RL further improving factual accuracy (Wang
83 et al. [2025b]).

84 Reliable RAG requires: (1) tightly coupled retrieval generation, (2) rigorous self-evaluation, and
85 (3) stable reward optimization. Prior approaches addressing these in isolation remain prone to
86 inconsistencies.

87 CREAM-RAG combines three core ideas: self-reflection-based adaptive retrieval, rank awareness
88 with training, and alignment between the retriever and generator. Together, these reduce extra retrieval
89 steps, make citations more accurate than using self-reflection alone, and give more reliable grounding
90 than recent consistency-based approaches.

91 **3 Methodology**

92 **3.1 Preprocessing and System Initialization**

93 Before extending the RAG framework, we perform key initialization steps: an instruction-tuned LLM
94 (LLaMA-2-7B) is used as both the generator and the self-reward model; a frozen copy of this model
95 is retained; and a vector database of external documents is incorporated for contextual retrieval.

96 **3.2 Consistency-Regularized Self-Rewarding Reinforcement Learning Within Training 97 Procedure**

98 We utilize CREAM to enhance the stability of RAG: unifying candidate generation, self-evaluation,
99 and regularization into a five-stage loop.

- 100 1. Retrieval: The RAG module retrieves relevant documents for the input prompt.
- 101 2. Candidate Generation: For the given query and retrieved context, the model produces
102 multiple responses conditioned on both prompt and documents.
- 103 3. Self-Scoring: The model ranks its own outputs, assigning preference-based rewards for
104 factual grounding, coherence, and relevance.
- 105 4. Consistency Regularization: A KL divergence loss is computed between the current model’s
106 preferences and those of a frozen reference model, penalizing inconsistent reward shifts to
107 stabilize learning.
- 108 5. Model Optimization: DPO updates the model parameters. The reference model is periodi-
109 cally updated to mirror the current model, maintaining training stability.

110 This end-to-end framework directly addresses key challenges of unstable rewards, retrieval noise, and
111 hallucinations by ensuring a stable training process aligned with high-quality output.

112 **3.3 Formalization of CREAM-RAG Objective Functions**

113 We now describe the mathematical underpinnings of CREAM-RAG’s training objectives, including
114 the reward function, ranking stability via Kendall Tau, and the final per-pair DPO loss.

115 First, the reward function r_{ij} reflects the improvement in logarithmic likelihood of the current model
116 P_θ over a reference model P_{ref} , with an optional normalization term.

$$r_{ij} = \beta \left[\log P_\theta(y_{ij} | x_j) - \log P_{\text{ref}}(y_{ij} | x_j) \right] + \beta \log Z(x_j) \quad (1)$$

117 To stabilize learning, we introduce a consistency measure between model and reference preferences
118 using Kendall Tau:

$$\tau_j = \frac{2}{N(N-1)} \sum_{1 \leq i < i' \leq N} \left[\begin{array}{l} \mathbb{1}((J_{ij} - J_{i'j})(K_{ij} - K_{i'j}) > 0) \\ -\mathbb{1}((J_{ij} - J_{i'j})(K_{ij} - K_{i'j}) < 0) \end{array} \right] \quad (2)$$

119 Finally, we train the model using a per-pair DPO loss, which aligns generation likelihood with
120 self-assessed preferences while preserving consistency with the reference model:

$$\begin{aligned}\mathcal{L}_{\text{DPO}}(\theta; y, y', x, z) = & -z(y, y', x) \log \left(\frac{P_\theta(y \mid x)}{P_{\text{ref}}(y \mid x)} \right) \\ & - (1 - z(y, y', x)) \log \left(\frac{P_\theta(y' \mid x)}{P_{\text{ref}}(y' \mid x)} \right)\end{aligned}\quad (3)$$

121 4 Experiments

122 4.1 Training and Evaluating

123 Our RAG contains 3 main parts for training and evaluation: long-form parsing, comprehension, and
124 hallucination training. In total, we used the full training and evaluation methods of 5 datasets (Hotpot-
125 QA, Natural Questions, Trivia-QA, RAGTruth, and SQuADv2)(Yang et al. [2018]; Kwiatkowski
126 et al. [2019]; Joshi et al. [2017]; Niu et al. [2024]; Rajpurkar et al. [2018]) and compared all scores
127 against a baseline. 10,000 samples were pulled from each dataset when training our baseline and
128 model. Additionally, when running tests on both our baseline and model, we utilized BERT scoring
129 to ensure effective semantic evaluation of both models (Zhang et al. [2020]). This allowed us to test
130 our model with strong, reliable metrics that resonate more with human judgment.

131 **Long-form parsing** tasks include two sets: a question-answer dataset that contains real user questions
132 from Google search where the answers are found in Wikipedia (Natural Questions; Kwiatkowski et al.
133 [2019]), and a multi-hop question-answer dataset with 113,000 Wikipedia-based question-answer
134 pairs (HotpotQA; Yang et al. [2018]). Both datasets used F1 and exact match as evaluation metrics;
135 we tested each dataset against a baseline.

136 **Comprehension tasks** include two question-answering datasets: TriviaQA (Joshi et al. [2017]) and
137 SQuAD v2 (Rajpurkar et al. [2018]). Trivia QA consists of over 650,000 trivia-style question-answer-
138 evidence triples, primarily to train the ability to answer factual knowledge-based questions. SQuAD
139 v2 contains over 50,000 unanswerable questions that look like answerable ones, challenging models
140 to both answer and filter between possible and impossible for questions. These two datasets gave the
141 same scoring metrics as HotpotQA and Natural Questions, Exact Match, and F1 score.

142 **Hallucination testing** was run on RAGTruth, which consists of more than 18,000 naturally generated
143 responses, annotated for evaluation of hallucinations (Niu et al. [2024]). RAGTruth uses similar data
144 evaluation metrics to the aforementioned datasets.

145 4.2 Baselines

146 We employed one main baseline to test against our CREAM-RAG model. Ablation testing of our
147 model occurred without the retrieval and CREAM.

148 Self-RAG (Asai et al. [2023]), is a self-rewarding language model that furthered Retrieval Augmented
149 Generation by adding a self-rewarding process. At first, RAG lacked versatility and struggled to
150 complete tasks without human-based reinforcement learning. Self-RAG was the first to break this
151 mold and depend upon Actor-Critic in RAG. We ran the same datasets for evaluation that we ran on
152 our model and compared the scoring.

153 5 Discussion

154 Our work introduces CREAM-RAG, a framework aimed at stabilizing reward signals and reducing
155 hallucinations. Our experiments across various long-form QA, comprehension, and hallucination-
156 specific tasks demonstrate how stabilizing the self-reward process leads to more reliable outputs, even
157 in noisy applications.

Table 1: Overall experimental results on four tasks. Balanced scores were calculated by allowing small token differences and numeric overlap. Token-based F1 with pre-processing took place, where the output text becomes normalized, allowing for minor capitalization or grammar differences. Normalized ground truth tokens are compared to normalized prediction tokens.

Benchmark	Balanced F1	Balanced EM	BERT Precision	BERT Recall	BERT F1
Llama-2-7B (No RAG)					
SQuADv2	34.8	39.4	24.2	24.6	24.4
NQ	32.9	32.6	34.2	35.1	27.6
TriviaQA	24.5	28.5	19.8	27.3	10.5
HotpotQA	21.1	30.5	24.0	16.3	10.9
RAGTruth	18.7	38.2	26.5	16.8	11.3
Self-RAG (RAG)					
SQuADv2	46.9	41.5	37.4	37.8	37.5
NQ	35.6	34.0	28.3	29.0	28.6
TriviaQA	26.3	31.4	44.9	41.9	43.3
HotpotQA	40.8	29.8	26.5	27.3	36.8
RAGTruth	56.3	36.5	11.5	28.5	20.1
CREAMRAG (RAG)					
SQuADv2	43.1	52.0	84.3	86.6	85.3
NQ	46.5	42.3	82.0	84.1	82.9
TriviaQA	45.2	41.7	83.5	82.9	83.1
HotpotQA	45.9	44.2	82.6	85.2	83.8
RAGTruth	82.7	43.0	84.5	86.2	85.3

158 5.1 Key Findings and Interpretation

159 Our results across multiple benchmark datasets show that CREAM-RAG improves answer quality in
 160 both short-answer and long-form question-answering tasks:

161 The model achieved high scores on TriviaQA, demonstrating accurate answer selection. On SQuAD
 162 v2, its performance shows a reduced tendency to hallucinate. For long-form questions on HotPotQA
 163 and Natural Questions, it maintained strong results, proving its ability to synthesize complex answers
 164 from multiple documents.

165 Our experimental results across a diverse set of data suggest that the key to improving factuality
 166 and reliability lies not only in better retrieval but in more stable reward modeling during training.
 167 When the reward signal is erratic, models struggle to learn consistent patterns of accuracy, leading to
 168 frequent hallucinations or degraded performance in noisy circumstances. Contrastingly, CREAM-
 169 RAG’s stabilized self-reward mechanism enables the model to better distinguish between accurate
 170 and inaccurate generations, even when retrieval results are partially irrelevant or convoluting.

171 These findings have several important implications. First, they underscore the critical role of reward
 172 signal quality in the success of RAG systems, specifically for tasks that require high factual accuracy.
 173 Second, they demonstrate that improving internal dynamics (e.g., stabilizing the self-reward process)
 174 can be as influential as external improvements such as improved retrieval or model scaling.

175 5.2 Broader Context and Significance

176 CREAM-RAG advances beyond traditional RAG by mitigating reward hacking and enables reliable
 177 self-evaluation, which is vital for creating autonomous RAGs in high-stakes domains.

178 Beyond its technical significance, CREAM-RAG paves the way for more autonomous and adaptive
 179 RAG systems. Its use of consistency regularization for self-evaluation allows deployment in fields
 180 with scarce human feedback, such as healthcare diagnostics, legal analysis, and scientific research.
 181 However, the self-rewarding functionality introduces serious ethical risks, including potential reward
 182 hacking, bias, and obscured accountability. To ensure safe and fair use, future applications require
 183 careful monitoring, clear reward logic, and domain-specific safeguards.

184 **5.3 Unexpected Observations**

185 Consistency regularization increased recall without degrading precision. The model retrieved more
186 relevant facts (e.g., 0.5980 recall on Natural Questions) without introducing noise (0.4147 precision),
187 indicating a new ability to take informed risks. Additionally, significantly higher balanced exact
188 match scores versus specific exact match confirm that answers are semantically, if not stringently,
189 correct.

190 **5.4 Metrics**

191 For scoring we used Balanced F1, Balanced EM, BERT Precision, BERT Recall, and BERT F1.
192 We utilized BERT score to give metrics on the models ability to perform in practical settings by
193 computing semantic similarity between the embeddings of predicted and reference sentences using
194 BERT models(utilized Roberta-Large)(Zhang et al. [2020]).

195 **5.5 Future Directions and Implications**

196 This paves the way for several applications: multimodal retrieval with unified self-evaluation across
197 text, images, and tables; reliable multi-step reasoning via chain-wide consistency checks; and hybrid
198 reward models that effectively blend automated scoring with limited human guidance. On top of
199 this, CREAM-RAG can be specialized for fields like medicine and law, increasing credibility by
200 mitigating hallucinations while also maintaining a strong comprehension of texts. This could allow
201 for AI trust and integration in important situations in a multitude of fields.

202 **6 Conclusion**

203 In this paper, we present CREAM-RAG, a consistency-regularized framework for Retrieval-
204 Augmented Generation that mitigates hallucinations and enhances factual reliability. By integrating
205 document retrieval with self-rewarding reinforcement learning and consistency loss, our approach sta-
206 bilizes training signals. Extensive experiments on long-form QA, comprehension, and hallucination
207 benchmarks demonstrate that CREAM-RAG delivers significant gains in output quality, particularly
208 when retrieval is imperfect. These results validate the role of consistency-based optimization in the
209 development of more reliable, autonomous models for high-stakes applications.

210 **7 Limitations**

211 Although CREAM-RAG sparks improvements in various benchmarks, our dependence on BERT
212 Score as a primary metric has limitations. Although BERT Score effectively captures semantic
213 similarity, it can overvalue the factuality of fluent or incorrect answers. This is especially problematic
214 in long-form or multi-hop QA (such as HotpotQA), where it may fail to identify hallucinations
215 or slight inaccuracies. Future evaluations would benefit from incorporating human judgment or
216 task-specific factual accuracy metrics to better gauge real-world applications.

217 Furthermore, while tested on five diverse datasets to demonstrate validity, CREAM-RAG is not
218 specialized for each, leading to performance variations. For instance, it achieves high precision on
219 TriviaQA but struggles with exact match scores on the multi-hop reasoning required by HotpotQA.
220 This diversity in datasets introduces variables like retrieval noise and answer styles, complicating
221 direct comparisons and contributing to inconsistent output. These results reinforce that, while
222 CREAM-RAG has broad applicability, achieving optimal performance on specialized tasks may
223 require specific tuning.

224 **References**

225 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to
226 retrieve, generate, and critique through self-reflection, 2023. URL <https://arxiv.org/abs/2310.11511>. arXiv preprint.

228 Shelly Bensal, Umar Jamil, Christopher Bryant, et al. Reflect, retry, reward: Self-improving llms via
229 reinforcement learning, 2025. URL <https://arxiv.org/abs/2505.24726>. arXiv preprint.

230 Mingyue Cheng, Yucong Luo, Jie Ouyang, Qi Liu, Huijie Liu, Li Li, Shuo Yu, Bohou Zhang, Jiawei
231 Cao, Jie Ma, Daoyu Wang, and Enhong Chen. A survey on knowledge-oriented retrieval-augmented
232 generation, 2025. URL <https://arxiv.org/abs/2503.10677>. arXiv preprint.

233 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Realm: Retrieval-
234 augmented language model pre-training, 2020. URL <https://arxiv.org/abs/2002.08909>.
235 arXiv preprint.

236 Hiroaki Hayashi, Prashant Budania, Peng Wang, Chris Ackerson, Raj Neervannan, and Graham
237 Neubig. Wikiasp: A dataset for multi-domain aspect-based summarization, 2020. URL <https://arxiv.org/abs/2011.07832>. arXiv preprint.

239 Gautier Izacard and Édouard Grave. Leveraging passage retrieval with generative models for open
240 domain question answering. In *Proceedings of EACL (Main Volume)*, 2021. URL <https://aclanthology.org/2021.eacl-main.74/>.

242 Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly
243 supervised challenge dataset for reading comprehension, 2017. URL <https://arxiv.org/abs/1705.03551>. arXiv preprint.

245 Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, et al. Training language models to self-correct via
246 reinforcement learning, 2024. URL <https://arxiv.org/abs/2409.12917>. arXiv preprint.

247 Tom Kwiatkowski, Jennimaria Palomaki, et al. Natural questions: A benchmark for question
248 answering research. *Transactions of the Association for Computational Linguistics*, 2019. URL <https://aclanthology.org/Q19-1026/>.

250 Youngwon Lee, Seung won Hwang, Daniel Campos, Filip Graliński, Zhewei Yao, and Yuxiong He.
251 Cord: Balancing consistency and rank distillation for robust retrieval-augmented generation, 2024.
252 URL <https://arxiv.org/html/2412.14581v1>. arXiv preprint.

253 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
254 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe
255 Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021. URL <https://arxiv.org/abs/2005.11401>. arXiv preprint arXiv:2005.11401.

257 Cheng Niu, Yuanhao Wu, Juno Zhu, et al. Ragtruth: A hallucination corpus for developing trustworthy
258 retrieval-augmented language models, 2024. URL <https://arxiv.org/abs/2401.00396>. arXiv
259 preprint.

260 Agada Joseph Oche, Ademola Glory Folashade, Tirthankar Ghosal, and Arpan Biswas. A systematic
261 review of key retrieval-augmented generation (rag) systems: Progress, gaps, and future directions,
262 2025. URL <https://arxiv.org/abs/2507.18910>. arXiv preprint.

263 Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don't know: Unanswerable questions
264 for squad, 2018. URL <https://arxiv.org/abs/1806.03822>. arXiv preprint.

265 Dan Su, Xiaoguang Li, Jindi Zhang, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. Read
266 before generate! faithful long form question answering with machine reading, 2022. URL
267 <https://arxiv.org/abs/2203.00343>. arXiv preprint.

268 Jinyu Wang, Jingjing Fu, Rui Wang, Lei Song, and Jiang Bian. Pike-rag: specialized knowledge
269 and rationale augmented generation, 2025a. URL <https://arxiv.org/abs/2501.11551>. arXiv
270 preprint.

271 Yuhao Wang, Ruiyang Ren, Yucheng Wang, Wayne Xin Zhao, Jing Liu, Huaqin Wu, and Haifeng
272 Wang. Reinforced informativeness optimization for long-form retrieval-augmented generation,
273 2025b. URL <https://api.semanticscholar.org/CorpusID:278910581>. Corpus record
274 (Semantic Scholar).

275 Zhaoyang Wang, Weilei He, Zhiyuan Liang, et al. Cream: Consistency regularized self-rewarding
276 language models, 2025c. URL <https://arxiv.org/abs/2410.12735>. arXiv preprint.

277 Zhilin Yang, Peng Qi, Saizheng Zhang, et al. Hotpotqa: A dataset for diverse, explainable multi-hop
278 question answering, 2018. URL <https://arxiv.org/abs/1809.09600>. arXiv preprint.

279 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and
280 Jason Weston. Self-rewarding language models, 2025. URL <https://arxiv.org/abs/2401.10020>. arXiv preprint.

282 Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav
283 Raja, Charlotte Zhuang, Dylan Slack, Qin Lyu, Sean Hendryx, Russell Kaplan, Michele Lunati,
284 and Summer Yue. A careful examination of large language model performance on grade school
285 arithmetic, 2024a. URL <https://arxiv.org/abs/2405.00332>.

286 Jinxu Zhang, Yongqi Yu, Yu Zhang, et al. Cream: Coarse-to-fine retrieval and multi-modal efficient
287 tuning for document vqa. In *Proceedings of the 32nd ACM International Conference on Multimedia*
288 (*MM '24*), 2024b. doi: 10.1145/3664647.3680750. URL <https://doi.org/10.1145/3664647.3680750>.

290 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating
291 text generation with bert, 2020. URL <https://arxiv.org/abs/1904.09675>.

292 Jiawei Zhou and Lei Chen. Openrag: Optimizing rag end-to-end via in-context retrieval learning,
293 2025. URL <https://arxiv.org/abs/2503.08398>. arXiv preprint.

294 Yiyang Zhou, Zhiyuan Fan, Dongjie Cheng, Sihan Yang, Zhaorun Chen, Chenhang Cui, Xiyao Wang,
295 Yun Li, Linjun Zhang, and Huaxiu Yao. Calibrated self-rewarding vision language models, 2024.
296 URL <https://arxiv.org/abs/2405.14622>. arXiv preprint.