

Revisiting Automated Evaluation for Long-form Table Question Answering

Yuqi Wang² Lyuhao Chen³ Songcheng Cai⁴ Zhijian Xu¹ Yilun Zhao¹

¹Yale University ²Independent Researcher ³Carnegie Mellon University ⁴University of Waterloo

Abstract

In the era of data-driven decision-making, Long-Form Table Question Answering (LFTQA) is essential for integrating structured data with complex reasoning. Despite recent advancements in Large Language Models (LLMs) for LFTQA, evaluating their effectiveness remains a significant challenge. We introduce LFTQA-Eval, a meta-evaluation dataset comprising 2,988 human-annotated examples, to rigorously assess the efficacy of current automated metrics in assessing LLM-based LFTQA systems, with a focus on faithfulness and comprehensiveness. Our findings reveal that existing automatic metrics poorly correlate with human judgments and fail to consistently differentiate between factually accurate responses and those that are coherent but factually incorrect. Additionally, our in-depth examination of the limitations associated with automated evaluation methods provides essential insights for the improvement of LFTQA automated evaluation.

 <https://github.com/YukiiWWW/LFTQA-Eval>

1 Introduction

In the current landscape where decisions are increasingly driven by data, the utility of tabular data provides a well-organized and efficient means of presenting data, which is essential for informed decision-making processes (Pasupat and Liang, 2015; Zhao et al., 2022a,b; Tang et al., 2024; Zhao et al., 2024a,b). Within this context, long-form table question answering (LFTQA) has emerged as a vibrant area of research, bridging the gap between structured data and the comprehensive insights required in real-world scenarios (Nan et al., 2022; Zhao et al., 2023a). As illustrated in Figure 1, given the complex question and numerous data points in a table, an LFTQA system must understand the relationships within the data and perform human-like

Title: United States House of Representatives Elections, 2012

District	Incumbent	First Elected	Party	Candidates
North Carolina 3	Walter Jones Jr	1994	Republican	Walter Jones Jr (R) 63.2% Erik Anderson (D) 36.8%
North Carolina 4	David Price	1996	Democratic	David Price (D) 74.4% Tim D'Annunzio (R) 25.6%
North Carolina 6	Howard Coble	1984	Republican	Howard Coble (R) 60.9% Tony Foriest (D) 39.1%
North Carolina 7	Mike McIntyre	1996	Democratic	Mike McIntyre (D) 50.1% David Rouzer (R) 49.9%
North Carolina 8	Larry Kissell	2008	Democratic	Richard Hudson (R) 54.1% Larry Kissell (D) 45.9%
North Carolina 10	Patrick McHenry	2004	Republican	Patrick McHenry (R) 57.0% Patsy Keever (D) 43.0%

 According to the voting result, which representative election in North Carolina districts was the most competitive, and why?

The race in the North Carolina 7th district was the most competitive, as the Democratic incumbent Mike McIntyre won by a slim margin, with only a 0.2% difference between him and his Republican challenger David Rouzer. Furthermore, this election was the only one among all North Carolina districts in 2012 that resulted in a margin of victory within less than 1 percent.



Figure 1: An example of the Long-form Table Question Answering (LFTQA) task investigated in our work.

reasoning over the tabular content to compose the paragraph-long answer.

Recent studies highlight the exceptional performance of Large Language Models (LLMs) in LFTQA tasks (Zhao et al., 2023b; Chen, 2023; Ye et al., 2023). However, the reliable evaluation of LLM-based systems in this domain remains a relatively unexplored area. Unlike conventional text generation tasks, where automatic metrics such as BLEU and ROUGE can somewhat effectively gauge the fluency and surface-level coherence of the generated text, LFTQA demands a more nuanced assessment approach. These traditional metrics, primarily designed for shorter texts, often fall short in LFTQA where the answers not only need to be contextually rich and structurally complex but also deeply rooted in logical reasoning derived from the underlying tabular data. They struggle to evaluate the logical structure and reasoning accuracy essential for long-form responses, as they do not account for the correctness of data interpretation or the ability to maintain a coherent argument over extended narratives. This limitation significantly impacts their utility in scenarios where the decision-making process relies heavily on the ac-

curate and logical processing of structured data, necessitating the development of new metrics that can more effectively measure these critical aspects.

Our research demonstrates that existing automatic metrics are inadequate in distinguishing between high-quality, factually accurate answers and those that are merely coherent. This discrepancy is alarming because developers might choose suboptimal systems for real-world applications if they rely solely on automatic metrics to compare and rank different LFTQA systems. To better investigate the automated evaluation methods for LFTQA tasks, we have constructed a meta-evaluation dataset named **LFTQA-Eval**, consisting of 2,988 human-annotated examples. Specifically, we gathered outputs from leading LLM-based systems on the FETAQA (Nan et al., 2022) and QTSUMM (Zhao et al., 2023a) datasets. We then benchmarked existing automatic evaluation metrics for these tasks, leveraging our collected human annotations across two distinct dimensions: faithfulness and comprehensiveness. Our experimental results demonstrate that all the examined automated metrics exhibit low correlations with human judgments, revealing their unreliability in determining the quality of LLM-generated answers and comparing different LLM-based systems. Moreover, we conducted an in-depth analysis of the failures associated with automated evaluation methods, supplemented by illustrative examples that provide valuable insights into potential areas for enhancement.

2 LFTQA-EVAL Construction

To better investigate the automated evaluation methods for LFTQA tasks, we have constructed a meta-evaluation dataset named LFTQA-Eval. The following subsections discuss the data collection methodology and annotation process.

2.1 Collecting LLM Output for LFTQA

We examine LFTQA automated evaluation methods on FETAQA and QTSUMM. Table 1 illustrates the basic data statistics of these two datasets.

- **FETAQA** (Nan et al., 2022) is designed for free-form table question answering, with answers averaging 18.9 words. It requires models to extract question-relevant information from the given table, and then aggregate and reason over this information to produce a coherent long-form answer.
- **QTSUMM** (Zhao et al., 2023a) requires models to perform reasoning and analysis akin to hu-

Property	FETAQA	QTSUMM
Table Source	Wikipedia	Wikipedia
Unique Tables	1,942	424
Avg. Rows per Table	14.2	12.0
Avg. Columns per Table	5.7	6.7
Avg. Table Title Length	5.4	7.4
Avg. Query Length	14.0	22.3
Avg. Summary Length	23.3	67.8
# QAs in Development Set	1,001	1,052

Table 1: Basic statistics of the FETAQA and QTSUMM test sets used in our experiments.

man thought processes on tables sourced from Wikipedia to produce paragraph-length answers. Compared to the FETAQA dataset, outputs in QTSUMM are longer, averaging 68.0 words.

Collecting LLM Output We collect output from eight popular LLMs using *zero-shot* prompting methods, including Llama-2&3-70B (Touvron et al., 2023), Qwen1.5-72B (Bai et al., 2023), Mistral-7B (Jiang et al., 2023a), DeepSeek-LLM-67B (DeepSeek, 2023), Gemma-7B (Team et al., 2024), OLMo-7B (Groeneveld et al., 2024), Yi-1.5-34B (01.AI, 2023), Phi-3-Medium (Abdin et al., 2024), and GPT-3.5-Turbo&4o (OpenAI, 2023). We use chat or instruct versions for each model. Additionally, we select the most recent, largest, and best-performing checkpoint available as of paper submission (i.e, June 15, 2024). We randomly sample 150 examples from the development sets of FETAQA and QTSUMM, and collect corresponding model outputs of these sampled examples. This results in a total of 2,988 examples within the final LFTQA-Eval benchmark (we exclude 12 empty model responses).

2.2 Evaluation Criteria

The automated evaluation of LFTQA tasks is challenging due to the unique features of LFTQA: 1) conducting intricate reasoning across multiple sources of information, and 2) ensuring factual accuracy while avoiding hallucination. To evaluate the reliability of automated evaluation methods for LFTQA, we collect human evaluation scores for each model output based on the the dimensions of **Faithfulness** and **Comprehensiveness**, respectively. Our preliminary study indicates that LLM-based systems exhibit the capability to generate texts that are both fluent and coherent, devoid of

spelling and grammatical errors. Therefore, we have excluded the evaluation of fluency and coherence from our analysis.

- **Faithfulness:** A good answer should be firmly rooted in the source table. It should consist of correct information from the table and precisely address the posed question, avoiding any inaccuracies or hallucinations.
- **Comprehensiveness:** A good answer should comprehensively include all relevant information from the tabular data that addresses the user’s question, fully meeting their informational needs.

2.3 Collecting Human Evaluation Scores

We tasked annotators to evaluate answers using a *Likert scale* ranging from 1 to 5 for the criteria of *faithfulness* and *comprehensiveness*, individually. To ensure the high quality of annotations, we hired eight undergraduate students proficient in English. Before starting the annotations, each annotator completed a one-hour online training session and reviewed a guide detailing the task execution steps. The annotators were compensated at an approximate hourly rate of \$10, aligned with the complexity and duration of the task. Each sample was independently evaluated by two different annotators to mitigate individual bias and variance in scoring. For each instance, we use the average of the two annotators’ scores as the final human evaluation score. Instances of significant disagreement (a variance greater than 2 points) were re-evaluated by an additional annotator. We achieved substantial inter-annotator agreements, with Krippendorff’s alpha for faithfulness- and comprehensiveness-level annotation at 0.738 and 0.714, respectively. This highlights the high-quality of LFTQA-Eval.

2.4 Collecting Automated Evaluation Scores

We examine following automatic metrics that are widely used in the LFTQA task, investigating their reliability in evaluating model performance:

- **BLEU** (Papineni et al., 2002) computes the geometric mean of the modified precision scores of the translated text and incorporates a brevity penalty factor. We use SacreBLEU (Post, 2018) for BLEU score calculation.
- **ROUGE** (Lin and Hovy, 2003) assesses the degree of lexical similarity between the generated text and the reference text. We employ F1 score for ROUGE-L.

- **METEOR** (Banerjee and Lavie, 2005) is developed to address the limitations of BLEU by introducing a method where alignment is established through the mapping of unigrams.
- **BERTScore** (Zhang et al., 2020) measures the similarity between the generated output and the reference text by utilizing contextualized token embeddings derived from a pre-trained model.
- **TAPAS-Acc** (Liu et al., 2022) assesses the faithfulness of table-to-text generation using TAPAS (Herzig et al., 2020) pretrained on the TabFact (Chen et al., 2020) dataset.
- **AutoACU** (Liu et al., 2023c) presents a reference-based automated evaluation system, utilizing atomic content units (ACUs) to gauge the similarity between text sequences.

We also adopt an LLM-based evaluation system, **G-Eval** (Liu et al., 2023a), to the LFTQA task. G-Eval employs LLMs using a chain-of-thought approach and the form-filling paradigm to assess the quality of generated text. We adopt the official prompt of evaluation instruction to assess the *faithfulness* and *comprehensiveness* of the generated answers, separately. The evaluation prompts used are presented in Appendix. We use the Llama-3.1-70B, GPT-4o-mini, and GPT-4o as the evaluators. For each model output collected in Section 2.1, we measure the scores of aforementioned metrics as automated evaluation scores.

2.5 Evaluating Automatic Evaluation Metrics

To evaluate the performance of automatic metrics, the human evaluation result on the same evaluation target is considered the gold standard, and metric performance is measured by the correlation between the human evaluation scores and automatic metric scores. Following previous work (Cohan and Goharian, 2016; Fabbri et al., 2021; Liu et al., 2023b), we calculate the correlation at the *instance-level*. Specifically, given n QA examples and m LFTQA systems, the human evaluation and an automatic metric result in two n -row, m -column score matrices H , M respectively. The instance-level correlation can be computed as the average of sample-wise correlations as follows:

$$r_{\text{ins}}(H, M) = \frac{\sum_i \mathcal{C}(H_i, M_i)}{n}, \quad (1)$$

Where H_i and M_i represent the evaluation results for the i -th data sample, with \mathcal{C} denoting a function that computes a correlation coefficient. In this

Metric	FETAQA		QTSUMM	
	Comp.	Faith.	Comp.	Faith.
BERT-Score	0.052	0.041	0.085	0.072
TAPAS-Acc	0.064	0.008	0.022	0.076
ROUGE	0.104	0.206	0.186	0.124
METEOR	0.191	0.216	0.204	0.187
AutoACU	0.143	0.262	0.268	0.188
BLEU	0.293	0.415	0.425	0.341
G-Eval _{4o-mini} Faith.	0.397	0.466	0.466	0.342
G-Eval _{4o-mini} Comp.	0.381	0.486	0.508	0.359
G-Eval _{3.1-70B} Comp.	0.424	0.521	0.543	0.426
G-Eval _{3.1-70B} Faith.	0.435	0.534	0.545	0.433
G-Eval _{4o} Faith.	0.470	0.628	0.633	0.470
G-Eval _{4o} Comp.	0.500	0.650	0.662	0.493

Table 2: Results of *instance-level* Pearson correlations between automatic metrics and human judgments on FETAQA and QTSUMM datasets.

study, we employ Pearson correlation to measure the correlations between human and automated evaluation systems.

3 Experimental Results

3.1 Main Results

Table 2 illustrates the instance-level Kendall’s tau correlation between automatic and human judgments. We can draw following two conclusions based on the results: **Existing automatic metrics fail in assessing the answers generated by LLM-based systems.** Table 2 reveals a general trend of low correlations across a range of metrics (e.g., BERT-Score, ROUGE, METEOR, and TAPAS-Acc), when evaluating individual LLM-generated responses. This indicates a widespread issue among current automatic metrics in measuring the faithfulness and comprehensive of LLM-generated answers, pointing to a systemic failure to align with human judgments at the instance level. **LLM-based metrics demonstrate a significant improvement over traditional automated metrics in terms of correlation with human evaluation.** G-Eval consistently achieves positive and high correlation scores, demonstrating the effectiveness of LLM-based metrics. Moreover, compared to Llama-3, GPT-4o yields higher scores, indicating that its evaluation results correspond more closely with human assessments. This superior performance reflects the enhanced evaluation capabilities of larger-size models in aligning with human judgment standards for the LFTQA task.

3.2 Case Study

To gain deeper insights into the failure cases of automated evaluation systems for LFTQA tasks, we conduct detailed human analyses by exploring scenarios where automated evaluations fall short. We identify three primary failure scenarios along with their underlying causes as follows:

The Effect of Question As we delve deeper into the examples, we observe that the clarity of the questions significantly impacts the quality of the generated answers. Ambiguous questions can lead the model to misinterpret the key elements, resulting in the retrieval of incorrect information from the tables. Furthermore, we discovered that some questions were subjective or open-ended, which led to a variety of perspectives and content in the answers. The information related to these questions may not be directly presented or elaborated in the given tables. Instead, it should be inferred and evaluated from external materials, requiring careful speculation and analysis. In contrast, both the ground truth and generated answers typically reflect only a subset of these potential viewpoints. Table 3 in Appendix presents detailed examples.

The Effect of Ground Truth Although ground truth is used as the standard reference in the evaluation process, it has certain issues that affect the quality of the assessment. Ground-truth answers often include extensive descriptive details, which can make them redundant and contain content irrelevant to the questions. Additionally, in some instances, the ground truth fails to provide the specific information requested in the question. This can lead to lower evaluation scores, even when the generated outputs are accurate. Table 4 in Appendix presents detailed examples.

The Effect of Generated Answer LLM-based models excel at incorporating additional, reasoning-intensive information that is not present in ground-truth answers. They generate direct, parallel structures in their responses, which align well with human expression in real-world applications. However, current automated metrics struggle to capture this supplementary information and concise structures, resulting in automated evaluation scores that are significantly lower than human scores. Table 5 in Appendix presents detailed examples.

4 Related Work

Table Question Answering (QA) challenges models to derive accurate answers from data presented in tables. This field is bifurcated into short-form QA (Pasupat and Liang, 2015; Zhong et al., 2017; Zhao et al., 2023c), which focuses on concise answers, and long-form QA (Nan et al., 2022; Zhao et al., 2023a), which demands the generation of paragraph-length responses that require advanced reasoning capabilities. Unlike short-form Table QA, where accuracy—the proportion of questions correctly answered—is the primary metric, the evaluation of LFTQA systems introduces unique challenges, particularly in ensuring the faithfulness and consistency of the generated responses to the underlying tabular data and posed questions. These responses must not only be accurate but also coherent and contextually relevant. In the era of LLMs (Zhao et al., 2023b; Chen, 2023; Ye et al., 2023), refining the automated evaluation methods to better capture the complexities of LFTQA remains a critical and ongoing area of research.

To evaluate automatic metric performance for text generation, several human evaluation benchmarks have been collected (Cohan and Goharian, 2016; Dhingra et al., 2019; Gabriel et al., 2021; Fabbri et al., 2021; Jiang et al., 2023b; Liu et al., 2024), comprising system-generated text and their human evaluation scores. The human evaluation result on the system-generated text is considered the gold standard, and metric performance is measured by the correlation between the human evaluation scores and automatic metric scores. To the best of our knowledge, we are the first to examine the automated evaluation methods for LFTQA research.

5 Conclusion

Our exploration into the evaluation of LLMs for LFTQA tasks reveals a significant gap between current automatic metrics and human judgment, particularly in assessing answer faithfulness and comprehensiveness. The insights from the LFTQA-Eval dataset highlight the need for more nuanced evaluation methods that align more closely with human evaluative standards. Addressing this discrepancy is essential for advancing the reliability of LFTQA systems and ensuring their practical utility in real-world scenarios.

Limitations

Our analysis is limited to 2,988 examples for which we have collected. While a larger evaluation set would allow for more statistically significant conclusions, this would require a significantly greater allocation of time and resources. We encourage future research to adopt our protocol and expand the benchmark for further analysis.

References

- 01.AI. 2023. [Yi: Open-source llm release](#).
- Marah Abidin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#).
- 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, Sheng-guang 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. *arXiv preprint arXiv:2309.16609*.
- 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.
- Wenhu Chen. 2023. **Large language models are few(1)-shot table reasoners**. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1120–1130, Dubrovnik, Croatia. Association for Computational Linguistics.
- Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyong Zhou, and William Yang Wang. 2020. **Tabfact: A large-scale dataset for table-based fact verification**. In *International Conference on Learning Representations*.
- Arman Cohan and Nazli Goharian. 2016. **Revisiting summarization evaluation for scientific articles**. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 806–813, Portorož, Slovenia. European Language Resources Association (ELRA).
- DeepSeek. 2023. Deepseek llm: Let there be answers. <https://github.com/deepseek-ai/DeepSeek-LLM>.
- Bhuvan Dhingra, Manaal Faruqui, Ankur Parikh, Ming-Wei Chang, Dipanjan Das, and William Cohen. 2019. **Handling divergent reference texts when evaluating table-to-text generation**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4884–4895, Florence, Italy. Association for Computational Linguistics.
- Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. **SummEval: Re-evaluating summarization evaluation**. *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Saadia Gabriel, Asli Celikyilmaz, Rahul Jha, Yejin Choi, and Jianfeng Gao. 2021. **GO FIGURE: A meta evaluation of factuality in summarization**. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 478–487, Online. Association for Computational Linguistics.
- Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muenighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, William Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah Smith, and Hannaneh Hajishirzi. 2024. **OLMo: Accelerating the science of language models**. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15789–15809, Bangkok, Thailand. Association for Computational Linguistics.
- Jonathan Herzig, Paweł Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. 2020. **TaPas: Weakly supervised table parsing via pre-training**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333, Online. Association for Computational Linguistics.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023a. **Mistral 7b**. *arXiv preprint arXiv:2310.06825*.
- Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao Huang, Bill Yuchen Lin, and Wenhu Chen. 2023b. **Tiger-score: Towards building explainable metric for all text generation tasks**.
- Chin-Yew Lin and Eduard Hovy. 2003. **Automatic evaluation of summaries using n-gram co-occurrence statistics**. In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 150–157.
- Ao Liu, Haoyu Dong, Naoaki Okazaki, Shi Han, and Dongmei Zhang. 2022. **PLOG: Table-to-logic pre-training for logical table-to-text generation**. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5531–5546, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruochen Xu, and Chenguang Zhu. 2023a. **G-eval: NLG evaluation using gpt-4 with better human alignment**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023b. **Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation**. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*

- (*Volume 1: Long Papers*), pages 4140–4170, Toronto, Canada. Association for Computational Linguistics.
- Yixin Liu, Alexander Fabbri, Jiawen Chen, Yilun Zhao, Simeng Han, Shafiq Joty, Pengfei Liu, Dragomir Radev, Chien-Sheng Wu, and Arman Cohan. 2024. [Benchmarking generation and evaluation capabilities of large language models for instruction controllable summarization](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4481–4501, Mexico City, Mexico. Association for Computational Linguistics.
- Yixin Liu, Alexander Fabbri, Yilun Zhao, Pengfei Liu, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023c. [Towards interpretable and efficient automatic reference-based summarization evaluation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16360–16368, Singapore. Association for Computational Linguistics.
- Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Victoria Lin, Neha Verma, Rui Zhang, Wojciech Kryściński, Hailey Schoelkopf, Riley Kong, Xiangru Tang, Mutethia Mutuma, Ben Rosand, Isabel Trindade, Renusree Bandaru, Jacob Cunningham, Caiming Xiong, Dragomir Radev, and Dragomir Radev. 2022. [FeTaQA: Free-form table question answering](#). *Transactions of the Association for Computational Linguistics*, 10:35–49.
- OpenAI. 2023. [Gpt-4 technical report](#). *ArXiv*, abs/2303.08774.
- 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.
- Panupong Pasupat and Percy Liang. 2015. [Compositional semantic parsing on semi-structured tables](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1470–1480, Beijing, China. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Xiangru Tang, Yiming Zong, Jason Phang, Yilun Zhao, Wangchunshu Zhou, Arman Cohan, and Mark Gestein. 2024. [Struc-bench: Are large language models good at generating complex structured tabular data?](#) 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 12–34, Mexico City, Mexico. Association for Computational Linguistics.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. [Gemma: Open models based on gemini research and technology](#).
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. 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 S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. 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, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu 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 fine-tuned chat models](#).

Yunhu Ye, Binyuan Hui, Min Yang, Binhua Li, Fei Huang, and Yongbin Li. 2023. [Large language models are versatile decomposers: Decomposing evidence and questions for table-based reasoning](#). In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '23*, page 174–184, New York, NY, USA. Association for Computing Machinery.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Yilun Zhao, Yunxiang Li, Chenying Li, and Rui Zhang. 2022a. [MultiHiertt: Numerical reasoning over multi hierarchical tabular and textual data](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6588–6600, Dublin, Ireland. Association for Computational Linguistics.

Yilun Zhao, Hongjun Liu, Yitao Long, Rui Zhang, Chen Zhao, and Arman Cohan. 2024a. [Financemath: Knowledge-intensive math reasoning in finance domains](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12841–12858, Bangkok, Thailand. Association for Computational Linguistics.

Yilun Zhao, Yitao Long, Hongjun Liu, Ryo Kamoi, Linyong Nan, Lyuhao Chen, Yixin Liu, Xiangru Tang, Rui Zhang, and Arman Cohan. 2024b. [DocMath-eval: Evaluating math reasoning capabilities of LLMs in understanding long and specialized documents](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16103–16120, Bangkok, Thailand. Association for Computational Linguistics.

Yilun Zhao, Linyong Nan, Zhenting Qi, Rui Zhang, and Dragomir Radev. 2022b. [ReasTAP: Injecting table reasoning skills during pre-training via synthetic reasoning examples](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9006–9018, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Yilun Zhao, Zhenting Qi, Linyong Nan, Boyu Mi, Yixin Liu, Weijin Zou, Simeng Han, Ruizhe Chen, Xiangru Tang, Yumo Xu, Dragomir Radev, and Arman Cohan. 2023a. [QTSumm: Query-focused summarization over tabular data](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1157–1172, Singapore. Association for Computational Linguistics.

Yilun Zhao, Haowei Zhang, Shengyun Si, Linyong Nan, Xiangru Tang, and Arman Cohan. 2023b. [Investigating table-to-text generation capabilities of large language models in real-world information seeking scenarios](#). In *Proceedings of the 2023 Conference on*

Empirical Methods in Natural Language Processing: Industry Track, pages 160–175, Singapore. Association for Computational Linguistics.

Yilun Zhao, Chen Zhao, Linyong Nan, Zhenting Qi, Wenlin Zhang, Xiangru Tang, Boyu Mi, and Dragomir Radev. 2023c. [RobuT: A systematic study of table QA robustness against human-annotated adversarial perturbations](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6064–6081, Toronto, Canada. Association for Computational Linguistics.

Victor Zhong, Caiming Xiong, and Richard Socher. 2017. [Seq2sql: Generating structured queries from natural language using reinforcement learning](#).

A Appendix

G-Eval for Evaluating Faithfulness

Task Introduction:

Given a complex question and a generated answer about a table, your task is to rate the answer's Faithfulness.

Evaluation Criteria:

Faithfulness(1-5): A good answer should accurately and completely address the given question. It must be based entirely on the information provided and should not include any unfaithful or hallucinated content.

Evaluation Steps:

1. Thoroughly review both the table and the question, ensuring a full understanding of the information they convey. Identify and analyze key points, critical data, and important details within the table that is relevant to the question.
2. Carefully examine the proposed answer, focusing on its faithfulness. Check for factual correctness and verify whether the answer reflects and aligns with the information presented in the table.
3. Evaluate the answer's faithfulness using a strict 1 to 5 rating scale, with 1 being the lowest and 5 the highest.

Figure 2: G-Eval for Evaluating the *faithfulness* of the LLM generated answer.

Error Type	Example	Explanation
Question is ambiguous	Question: Who were the top three scorers for the 1961-62 Michigan Wolverines men's basketball team and how many points did they score?	It may take individual scores but is phrased in a way that could be interpreted as asking for a total score, potentially leading to the total score being treated as another player in the ranking.
Subjective issues	Question: How did the performance of Tom Brady in terms of passing yards during the Regular Season 2011 compare with other quarterbacks listed in 2011?	The subject of these questions might result in multiple reasonable interpretations and answers. For example, responses could pertain to Tom's scoring rate, passing rate, ball handling performance, etc., each in different ways.
Open-ended questions	Question 1: Summarize the basic information of the episode(s) written by Damon Lindelof. Question 2: Summarize the performance of Weekend Hussler in the Caulfield Guineas.	These questions involve various perspectives and require external information to be adequately addressed. For example, the first question might pertain to understanding the play, including plot trends, character development, and thematic elements in the episode. Different background knowledge and perspectives will result in varying answers.

Table 3: Case studies on evaluation errors due to the effects of questions.

Error Type	Example	Explanation
Redundancy in Ground Truth Answer	Generated Answer: The quantity of the aircraft models that were first introduced between the years 1980 and 1985 inclusive is 3 and their build years are 1978, 1979, and 1971. Ground Truth Answer: The number of aircraft models that were initially launched within the timeframe extending from the year 1980 up to and including the year 1985 amounts to a total of three distinct models. Specifically, the construction years for these models are identified as the years 1978, 1979, and 1971 respectively.	Ground truth answers tend to include additional details that may not directly contribute to answering the question in some cases while the generated answers are concise with all essential information. Based on the example provided, the ground truth answer offers a more comprehensive and detailed description of the launch times of the aircraft models.
Mismatch Between Question and Ground Truth Answer	Generated Answer: The track that has the lowest bpm on the Speed SFX series is 'Rhythm & Police (K.O.G G3 Mix)' and the track with the highest bpm on the Speed SFX series is the 'Mission: Impossible Theme'. Ground Truth Answer: The Speed SFX's lowest bpm is 175 on tracks 5, and the highest bpm is 195 on 20 in the Speed series.	The ground truth focuses on the BPM values without mentioning specific track names, which is the primary detail the question asks for. Conversely, the prediction accurately identifies the tracks by name, directly addressing the question. Therefore, the ground truth fails to provide the specific information that the question seeks.

Table 4: Case studies on evaluation errors due to the effects of ground truth answers.

Error Type	Example	Explanation
Inconsistency in LLM Output Length	<p>Generated Answer: The HP LaserJet 4250 has a faster CPU speed of 460 MHz compared to the 4200's 300 MHz and its print speed is also higher at 43 pages per minute (PPM) compared to the 4200's 35 PPM, yet both printers come with a standard memory of 48 MB and can be expanded up to a maximum of 512 MB.</p> <p>Ground Truth Answer:The HP LaserJet 4200 and 4250 models, both in HP LaserJet 4000 series, have some same and different specifications. In same parts, both models come out in November 2002 and stop in June/November 20051. They have print resolution of 1200 DPI. Also, both printers can make the memory bigger, up to 512 MB. But, there are some different things between these two models. The LaserJet 4200 has 300 MHz CPU speed, but 4250 has faster 460 MHz CPU. Also, 4200 model can print 35 pages in one minute (PPM), but 4250 can print a little faster, 43 PPM".</p>	<p>The GT length is 68, while the LLM-generated length is 52.16. This indicates that the predicted output is significantly shorter than the actual answer. Generated responses tend to be concise and straightforward, focusing on delivering key points efficiently. In contrast, the actual answer provides more extensive information, with greater detail and elaboration. This difference highlights a tendency for automated responses to prioritize brevity.</p>
Answers' Structures	<p>Generated Answer: The quantity of the aircraft models that were first introduced between the years 1980 and 1985 inclusive is 3 and their build years are 1978, 1979, and 1971.</p> <p>Ground Truth Answer: Between the years 1980 to 1985 altogether, Agderfly added three airplane models to its fleet. In the year 1980, one Piper Chieftain made in 1978 was added, also one Piper Tomahawk was made in 1979 in the same year. The 1985 year, one Piper Seneca which was made in 1971. In total, during this time, Agderfly added three aircraft models whose combined quantity is four units.</p>	<p>Generated answers tend to be structured with parallel objects, while ground truth answers often utilize complex clauses to introduce related information thoroughly. In this example, the generated answer simply lists the years, while the ground truth introduces the information for each year in a single, comprehensive sentence. This discrepancy in structure can result in misalignment between automated predictions and the expected answers, impacting the accuracy of evaluations and interpretations.</p>

Table 5: Case studies on evaluation errors due to the effects of generated answers.

G-Eval for Evaluating Comprehensiveness

Task Introduction:
Given a complex question and a generated answer about a table, your task is to rate the answer's Comprehensiveness.

Evaluation Criteria:
Comprehensiveness(1-5): A good answer should provide all the necessary information to address the question comprehensively. Additionally, it should avoid including details that, while consistent with the tabular data, are irrelevant to the given question.

Evaluation Steps:

1. Carefully review the table and the question, ensuring you understand the full scope of the information provided. Identify all relevant points and details necessary to answer the question comprehensively.
2. Analyze the proposed answer to determine if it covers all the key aspects and addresses the question fully. Check whether the answer omits any important information or includes unnecessary details.
3. Evaluate the answer's comprehensiveness using a 1 to 5 rating scale, where 1 indicates the least comprehensive and 5 indicates the most.

Figure 3: G-Eval for Evaluating the *Comprehensiveness* of the LLM generated answer.