

- Li Du, Xiao Ding, Yue Zhang, Ting Liu, and Bing Qin. 2022. [A graph enhanced BERT model for event prediction](#). In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 2628–2638. Association for Computational Linguistics.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. [Gptscore: Evaluate as you desire](#). *CoRR*, abs/2302.04166.
- Alberto García-Durán, Sebastijan Dumancic, and Mathias Niepert. 2018. [Learning sequence encoders for temporal knowledge graph completion](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 4816–4821. Association for Computational Linguistics.
- Mark Granroth-Wilding and Stephen Clark. 2016. [What happens next? event prediction using a compositional neural network model](#). In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*, pages 2727–2733. AAAI Press.
- Danny Halawi, Fred Zhang, Chen Yueh-Han, and Jacob Steinhardt. 2024. [Approaching human-level forecasting with language models](#). *CoRR*, abs/2402.18563.
- Woojeong Jin, Rahul Khanna, Suji Kim, Dong-Ho Lee, Fred Morstatter, Aram Galstyan, and Xiang Ren. 2021. [Forecastqa: A question answering challenge for event forecasting with temporal text data](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 4636–4650. Association for Computational Linguistics.
- Ehsan Kamalloo, Nouha Dziri, Charles L. A. Clarke, and Davood Rafiei. 2023. [Evaluating open-domain question answering in the era of large language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 5591–5606. Association for Computational Linguistics.
- Tom Kocmi and Christian Federmann. 2023. [Large language models are state-of-the-art evaluators of translation quality](#). In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation, EAMT 2023, Tampere, Finland, 12-15 June 2023*, pages 193–203. European Association for Machine Translation.
- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Manling Li, Sha Li, Zhenhailong Wang, Lifu Huang, Kyunghyun Cho, Heng Ji, Jiawei Han, and Clare R. Voss. 2021a. [The future is not one-dimensional: Complex event schema induction by graph modeling for event prediction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 5203–5215. Association for Computational Linguistics.
- Zhongyang Li, Xiao Ding, and Ting Liu. 2018. [Constructing narrative event evolutionary graph for script event prediction](#). In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pages 4201–4207. ijcai.org.
- Zixuan Li, Xiaolong Jin, Wei Li, Saiping Guan, Jiafeng Guo, Huawei Shen, Yuanzhuo Wang, and Xueqi Cheng. 2021b. [Temporal knowledge graph reasoning based on evolutionary representation learning](#). In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, pages 408–417. ACM.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. [Let’s verify step by step](#). *CoRR*, abs/2305.20050.
- Sijia Liu, Patrick Lange, Behnam Hedayatnia, Alexandros Papangelis, Di Jin, Andrew Wirth, Yang Liu, and Dilek Hakkani-Tur. 2023a. [Towards credible human evaluation of open-domain dialog systems using interactive setup](#). In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 13264–13272. AAAI Press.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 2511–2522. Association for Computational Linguistics.
- Yixin Liu, Alexander R. Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023c. [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)*, ACL 2023, Toronto, Canada, July 9-14, 2023, pages 4140–4170. Association for Computational Linguistics.
- Ruilin Luo, Tianle Gu, Haoling Li, Junzhe Li, Zicheng Lin, Jiayi Li, and Yujiu Yang. 2024. [Chain of history: Learning and forecasting with llms for temporal knowledge graph completion](#). *CoRR*, abs/2401.06072.
- Yunshan Ma, Chenchen Ye, Zijian Wu, Xiang Wang, Yixin Cao, Liang Pang, and Tat-Seng Chua. 2023. [Structured, complex and time-complete temporal event forecasting](#). *CoRR*, abs/2312.01052.
- 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 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 12076–12100. Association for Computational Linguistics.
- Xiao Pu, Mingqi Gao, and Xiaojun Wan. 2024. [Is summary useful or not? an extrinsic human evaluation of text summaries on downstream tasks](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 9389–9404. ELRA and ICCL.
- Fengcai Qiao, Pei Li, Jingsheng Deng, Zhaoyun Ding, and Hui Wang. 2015. [Graph-based method for detecting occupy protest events using GDELT dataset](#). In *2015 International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery, CyberC 2015, Xi'an, China, September 17-19, 2015*, pages 164–168. IEEE Computer Society.
- Jie Ruan, Xiao Pu, Mingqi Gao, Xiaojun Wan, and Yuesheng Zhu. 2024. [Better than random: Reliable NLG human evaluation with constrained active sampling](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 18915–18923. AAAI Press.
- Xiaoming Shi, Siqiao Xue, Kangrui Wang, Fan Zhou, James Zhang, Jun Zhou, Chenhao Tan, and Hongyuan Mei. 2023. [Language models can improve event prediction by few-shot abductive reasoning](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Lihong Wang, Juwei Yue, Shu Guo, Jiawei Sheng, Qianren Mao, Zhenyu Chen, Shenghai Zhong, and Chen Li. 2021. [Multi-level connection enhanced representation learning for script event prediction](#). In *WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021*, pages 3524–3533. ACM / IW3C2.
- Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang, Cunxiang Wang, Hao Chen, Chaoya Jiang, Rui Xie, Jindong Wang, Xing Xie, Wei Ye, Shikun Zhang, and Yue Zhang. 2023. [Pandalm: An automatic evaluation benchmark for LLM instruction tuning optimization](#). *CoRR*, abs/2306.05087.
- Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024. [Back to the future: Towards explainable temporal reasoning with large language models](#). In *Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024*, pages 1963–1974. ACM.
- Mengqi Zhang, Yuwei Xia, Qiang Liu, Shu Wu, and Liang Wang. 2023. [Learning long- and short-term representations for temporal knowledge graph reasoning](#). In *Proceedings of the ACM Web Conference 2023, WWW 2023, Austin, TX, USA, 30 April 2023 - 4 May 2023*, pages 2412–2422. ACM.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Liang Zhao. 2022. [Event prediction in the big data era: A systematic survey](#). *ACM Comput. Surv.*, 54(5):94:1–94:37.
- Pengpeng Zhou, Bin Wu, Caiyong Wang, Hao Peng, Juwei Yue, and Song Xiao. 2022. [What happens next? combining enhanced multilevel script learning and dual fusion strategies for script event prediction](#). *Int. J. Intell. Syst.*, 37(11):10001–10040.
- Fangqi Zhu, Jun Gao, Changlong Yu, Wei Wang, Chen Xu, Xin Mu, Min Yang, and Ruifeng Xu. 2023. [A generative approach for script event prediction via contrastive fine-tuning](#). In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 14056–14064. AAAI Press.
- Andy Zou, Tristan Xiao, Ryan Jia, Joe Kwon, Mantas Mazeika, Richard Li, Dawn Song, Jacob Steinhardt, Owain Evans, and Dan Hendrycks. 2022. [Forecasting future world events with neural networks](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.

	Train	Dev	Test
CEs	41,424	474	1,519
AQA	82,846	948	3,038
STF	54,619	469	1,588
LTF	72,355	948	3,038
MCNC	86,326	886	3,697
MCAC	82,846	948	3,038
VQA	86,326	886	3,697

Table 4: Statistics of the data splitting of OpenForecast.

A Dataset

A.1 Details of Dataset Collection

In this work, the Wikipedia dump (20240320) and WCEP data before 2024 April are selected as the original sources. For the Wikipedia dump, we initially get 4,056,152 articles after excluding those with fewer than 200 words. We then implement a multi-step filtration method to efficiently filter event-related articles. The steps are as follows: (1) **Section Title Filtering.** A statistical analysis of section titles across Wikipedia articles reveals that event-related articles often contain sections such as "Background", "Development", "Aftermath", and "Reaction". Consequently, we collect all the section titles that may be relevant to the event and filter articles containing these section titles. (2) **Category Filtering.** While the section title filtering step excludes most non-event articles, a significant amount of noisy articles remained. To address this, we employ the Mixtral-8x7b to categorize the articles, discarding those that belong to individuals, locations, organizations, nations, etc. (3) **Time Filtering.** Events from different periods exhibit distinct evolution patterns. So we leverage Mixtral-8x7b to identify the occurrence dates of events and filter out those before 1950. For the WCEP dataset, which only documents current events after 2000, filtration is not employed. We first crawl and extract all metadata including event summaries, external links, time of occurrence, event categories, and subheaders from the WCEP website. Then we leverage the Newspaper3k project to scrape the news articles of external links. Finally, events sharing the same subheaders are aggregated into complex events.

By category filtering in data collection, we recognize the regular events (such as annual festivals, exhibitions, and conferences) by LLMs and the number is only 925/43419 \approx 2.1%. Coupled with the diversity of problems, the impact of regular

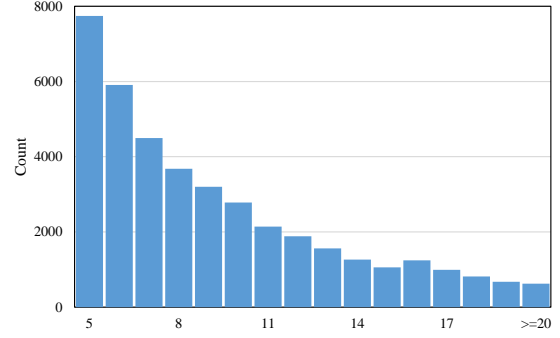


Figure 8: Number of complex events with varying atomic events number.

events is very small.

The prompt templates for the dataset construction are shown in Table 11.

A.2 Details of Dataset Splitting

To prevent knowledge leakage for the evaluation, we take the data before 2023/06/30 as the trainset, data between 2023/07/01 and 2023/08/31 as the validation set, and data between 2023/09/01 and 2024/03/31 as the testset. The detailed statistics of data splits of OpenForecast are presented in Table 4. Figure 8 illustrates the distribution of atomic event counts of complex events. Long-term complex events typically encompass more atomic events.

A.3 The Quality of the Event Timeline Construction

To evaluate the event timeline construction pipeline, which includes dataset collection and event timeline annotation, we randomly select 200 complex events (the selected samples encompass all event types with various article lengths) corresponding to 200 event timelines and ask two human annotators⁶ to evaluate the extracted event timelines from multiple dimensions, as outlined below:

- *Event Relevance* examines whether non-event articles, including those related to individuals, locations, organizations, nations, festivals, entertainment, concepts, etc from Wikipedia, are successfully excluded.
- *Completeness* evaluates the completeness of the extracted event timelines. Annotators first review the articles to grasp the overall event timeline and then determine whether the extracted timelines

⁶Our annotation team consists of two graduate students engaged in information processing. Another annotator will review their annotation results and eliminate their discrepancies.