

Foundations of Global Consistency Checking with Noisy LLM Oracles

Anonymous ACL submission

Abstract

Ensuring that collections of natural-language facts are globally consistent is essential for tasks such as fact-checking, summarization, and knowledge base construction. While Large Language Models (LLMs) can assess the consistency of small subsets of facts, their judgments are noisy, and pairwise checks are insufficient to guarantee global coherence. We formalize this problem and show that verifying global consistency requires exponentially many oracle queries in the worst case. To make the task practical, we propose an adaptive divide-and-conquer algorithm that identifies minimal inconsistent subsets (MUSes) of facts and optionally computes minimal repairs through hitting-sets. Our approach has low-degree polynomial query complexity. Experiments with both synthetic and real LLM oracles show that our method efficiently detects and localizes inconsistencies, offering a scalable framework for linguistic consistency verification with LLM-based evaluators.

1 Introduction

Ensuring the *global consistency* of sets of natural-language facts is essential for core NLP applications such as multi-document summarization, fact-checking, and knowledge base construction (Chen et al., 2024; Guo et al., 2022). For example, reports describing the same real-world event may contain overlapping or partially conflicting claims; systems must determine whether all claims can jointly hold, and if not, identify where contradictions arise. Crucially, it is not enough to merely detect that some inconsistency exists. In many downstream pipelines, naively discarding all facts whenever a contradiction is detected is unacceptable: a single conflict can cause large numbers of otherwise correct statements to be removed, degrading the quality of summaries, reports, or databases. What is needed instead are *explanations of inconsistency* and *principled ways to repair* fact sets by retaining as many mutually consistent facts as possible while isolating the smallest conflicting groups.

This motivates a focus on *Minimal Unsatisfiable Subsets* (MUSes), the smallest sets of claims that cannot jointly be true.

Large Language Models (LLMs) are increasingly used as *judges* for evaluation and verification tasks (Gu et al., 2025; Zhu et al., 2025; Wang et al., 2024a), and they can often assess whether a *small* set of claims appears consistent (Hong et al., 2025; Li et al., 2024). However, pairwise checks do not imply global consistency, and exhaustively querying all subsets is infeasible (Kumar et al., 2023). Moreover, direct “all-at-once” judgments become increasingly unreliable as the number of claims grows, due to longer inputs and denser interactions among statements. The central challenge is therefore: *how can we verify global consistency over many claims while issuing as few noisy LLM-judge calls as possible, and while retaining fine-grained explanations of inconsistency?*

Prior NLP work on factuality, contradiction detection, and fact verification has largely operated at the level of individual claims or pairs—for example, verifying a claim against evidence (Wang and Shu, 2023; Tan et al., 2025), decomposing complex claims (Pan et al., 2023), or retrieving supporting passages (Aly and Vlachos, 2022; de Marneffe et al., 2008). Classic surveys emphasize the importance of factual consistency in NLP systems (Thorne and Vlachos, 2018). However, these approaches do not address *global, multi-statement inconsistency detection*: ensuring that a set of extracted facts is jointly coherent. In real NLP pipelines—such as multi-document retrieval-augmented generation, large-scale information extraction, or report generation—the extracted fact set itself can become contradictory even when each claim is individually supported. From a computational perspective, global consistency is intractable in general (Lee and Leung, 2010), but the presence of small conflict sets or structural regularities makes adaptive, query-efficient approaches viable in practice.

In this paper, we formalize scalable global con-

sistency verification as querying a *noisy subset-consistency oracle*, instantiated by an LLM judge. We show that global consistency cannot be certified from pairwise checks alone and that worst-case query complexity is exponential even under strong assumptions. To make the task practical, we propose an adaptive divide-and-conquer algorithm that localizes Minimal Unsatisfiable Subsets (MUSes) within a set of natural-language claims and, when desired, computes minimal repairs via hitting-set duality. We provide theoretical bounds on query complexity and noise amplification, and empirically demonstrate that the proposed method efficiently detects and explains contradictions using LLMs, substantially improving recall over direct “all-at-once” judging while preserving high precision.

2 Related Work

Local Consistency Verification. Fact and claim consistency verification with LLMs has attracted growing interest, motivated by challenges such as hallucinations and misinformation (Rahman et al., 2025; Singhal et al., 2024). Early methods focused on extracting structured knowledge units such as subject–predicate–object triples from both LLM outputs and reference texts to detect local inaccuracies and support retrieval-augmented verification (Chen et al., 2025; Cao et al., 2025; Lewis et al., 2021). Other approaches decompose responses into atomic claims that are checked independently, which works well for short texts but struggles in long-form or multi-document settings (Hu et al., 2025; Wanner et al., 2024). Reranking methods (Liu et al., 2025) can help mitigate this limitation, but typically require access to model internals, limiting their practical deployment.

Global Consistency and Contradictions. Beyond knowledge-unit extraction, much work has studied factuality and contradiction detection, often casting the problem as a natural-language inference (NLI) task (Thorne et al., 2018; Kryscinski et al., 2020). Benchmarks such as FEVER and VitaminC emphasize identifying local entailment or contradiction, but they do not address whether an entire set of claims can jointly be true (Thorne et al., 2018; Schuster et al., 2021). More recently, LLMs themselves have been used as judges for factuality and coherence (Zheng et al., 2023; Wang et al., 2024b), extending this line of work beyond classifier-based approaches. Concurrent work has studied logical consistency of LLMs on propositional queries over knowledge graphs (Ghosh et al.,

2025), e.g. whether $A \wedge B$ is judged consistent with A and B separately, improving performance via fine-tuning.

To our knowledge, our work is the first to formalize *scalable global consistency verification*, prove theoretical query-complexity bounds, and propose adaptive algorithms for isolating minimal inconsistent subsets under noisy LLM oracles. We view this as the first step toward principled, scalable methods for global consistency verification with LLMs.

3 Setting the Stage

Problem Definition. Given a finite fact set $F = \{f_1, \dots, f_N\}$, we assume that there exists a *ground truth function* $A : 2^F \rightarrow \{\text{cons}, \text{incons}\}$, i.e., takes a (sub)set of facts and returns whether it is globally consistent (cons) or not (incons).

Let F be a finite set of facts and $C = \{C_1, \dots, C_m\}$ a family of scopes with $C_i \subseteq F$. We seek a kept set $F' \subseteq F$ that maximizes coverage while satisfying all per-scope constraints:

$$\max_{F' \subseteq F} |F'| \quad \text{s.t.} \quad A(F' \cap C_i) = \text{cons}, \quad \forall i \in [m].$$

Define $\tilde{C}_i \triangleq F' \cap C_i$. Then $\tilde{C}_i \subseteq C_i$ and $A(\tilde{C}_i) = \text{cons}$ for all i , and $|\bigcup_{i=1}^m \tilde{C}_i| = |F'|$. An equivalent formulation of the above objective in terms of minimizing a size of *hitting set* is formulated in Appendix A.1.¹

Complexity Landscape. In the worst case, solving the above optimization problem is NP-hard, as A can be used to encode any function (such as boolean satisfiability). Note that it is also possible for all fact pairs to be mutually consistent while the full set remains globally inconsistent (see Appendix A.2).

LLM as Noisy Subset-Consistency Oracle. We simulate A using LLMs. Given a finite fact set $F = \{f_1, \dots, f_N\}$, we model a pretrained LLM as a *noisy subset-consistency oracle* (NSCLM) O as follows. For any subset $S \subseteq F$ we prompt “Are the following claims mutually consistent?” and receive a stochastic response $O(S) \in \{\text{cons}, \text{incons}\}$. Let α, β denote the error rates on the oracle O ’s performance.

$$\Pr[O(S) = \text{incons} \mid A(S) = \text{cons}] \leq \alpha$$

$$\Pr[O(S) = \text{cons} \mid A(S) = \text{incons}] \leq \beta.$$

When using LLM as O , naively querying $O(F)$ is unreliable in practice as the set size of F increases,

¹The hitting set problem seeks the smallest subset of elements that intersects every set in a given family.

and as noted above pairwise checks are insufficient to detect inconsistencies.

4 Method

We now formalize our approach to consistency checking under a NSCLM. Our algorithm assumes as input a set of natural-language facts F and a family of *constraints* C with $C_j \subseteq F$ for all $C_j \in C$. Constraints may be given externally (e.g., from a schema or ontology) or constructed automatically from F ; here we focus on the given-constraint case for clarity. We start with the definition of Minimal Unsatisfiable Subset (MUS) which forms the basic building block of our approach.

Definition 4.1 (Minimal Unsatisfiable Subset w.r.t. Oracle O). A subset $U \subseteq F$ is an Minimal Unsatisfiable Subset (MUS) if $O(U) = \text{incons}$ but $O(U') = \text{cons}$ for all proper subsets $U' \subset U$.

Our procedure (Algorithm 1) runs an iterative two-step loop—*MUS extraction* followed by *greedy repair via a hitting set*—repeating until all constraints are consistent, and return the surviving facts F' . The soundness guarantee of the procedure is presented in Appendix A.3.

Assumptions. Our theoretical analysis assumes (i) approximate independence of repeated oracle calls so that majority voting reduces noise, (ii) small conflict size k in practice (typically $k \leq 3$), and (iii) constraint scopes C_i that are either externally defined or automatically constructed from entity or event clusters. These assumptions are used *only* to derive worst-case guarantees and are *not required* by the empirical method: all experiments use a single oracle call per query ($r = 1$) with automatically constructed scopes.

MUS Extraction via Divide-and-Conquer. Our MUS localization procedure builds on the QuickXplain algorithm (Junker, 2004), which given a set of possibly inconsistent constraints, identifies a minimal unsatisfiable subset through a recursive divide-and-conquer strategy. By recursively partitioning the constraint set and reusing intermediate results, QuickXplain achieves logarithmic query growth in subset size. We discuss more details about QuickXplain in Appendix C.

Given an input (O, S, B) , where B denotes a background set of facts assumed to be consistent, QuickXplain begins by splitting S into two parts, $S_1 \cup S_2$, and querying the oracle O on $B \cup S_1$. If $O(B \cup S_1) = \text{incons}$, it continues by recursing on $(S_1, B \cup S_2)$; otherwise, it proceeds with $(S_2, B \cup S_1)$. The recursion keeps narrowing the search until

Algorithm 1: QXR

Input: Facts F ;

Constraints $C = \{C_1, \dots, C_m\}$;

Noisy LLM oracle O

Output: Consistent facts F'

```

1  $F' \leftarrow F$ ;
2  $\mathcal{O}_{\text{incons}} \leftarrow \emptyset$ ;
3 while  $\exists C_j \in C : O(C_j) = \text{incons}$  do
4    $\mathcal{U} \leftarrow \emptyset$ ;
5   for  $C_j \in C$  with  $O(C_j) = \text{incons}$  do
6      $\mathcal{U} \leftarrow \mathcal{U} \cup \{\text{QX}(O, C_j, \emptyset)\}$ 
7    $H \leftarrow \text{GREEDYHITTINGSET}(\mathcal{U})$ ;
8    $F' \leftarrow F' \setminus H$ ;
9    $C \leftarrow \{C_j \setminus H : C_j \in C\}$ ;
10   $\mathcal{O}_{\text{incons}} \leftarrow \mathcal{O}_{\text{incons}} \cup \mathcal{U}$ 
11 return  $F'$ 

```

$|S| = 1$, ultimately returning a subset-minimal inconsistent set U .

Assume S contains a MUS U of size k . Starting from $\text{QX}(O, S, \emptyset)$, the QuickXplain procedure returns some MUS $U' \subseteq S$ using at most $\mathcal{O}(k \log |S|)$ oracle calls to O .

Greedy Repair via Hitting Set (Minimal Correction Set). Once MUSes are extracted, we identify the inconsistent scopes:

$$T_{\text{incons}} = \{j \in [m] : O(C_j) = \text{incons}\}. \quad (1)$$

For each $j \in T_{\text{incons}}$, we obtain a MUS $U_j \subseteq C_j$, and form the family of conflicts $\mathcal{U} = \{U_j : j \in T_{\text{incons}}\}$. We then compute a *repair set* $H \subseteq F$ that intersects every MUS:

$$\forall U \in \mathcal{U}, \quad H \cap U \neq \emptyset. \quad (2)$$

Such a minimal hitting set H corresponds exactly to a *Minimal Correction Set (MCS)*—the smallest subset of facts whose removal restores global consistency. We remove H to obtain the consistent subset $F' = F \setminus H$. Intuitively, the hitting set selects the fewest facts that “break” all discovered inconsistencies. For example, if $\mathcal{U} = \{\{a, b, c\}, \{a, d, e\}\}$, then any H intersecting both conflicts is valid, and the minimal hitting set $H = \{a\}$ yields the maximal consistent subset $F' = \{b, c, d, e\}$. In practice we use a greedy solver that iteratively selects the fact covering the largest number of uncovered MUSes, achieving the optimal logarithmic approximation ratio for this NP-hard problem but can be approximated efficiently (see Appendix B).

Let N be number of facts in F , $m = |C|$ be the number of constraints, k the maximum size of

Model	VitaminC						FEVER					
	Direct (baseline)			QXR (ours)			Direct (baseline)			QXR (ours)		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Claude 3.7 (Anthropic, 2025a)	0.979	0.854	0.909	0.956	0.975	0.965	0.992	0.805	0.873	0.983	0.977	0.980
Claude 4 (Anthropic, 2025b)	0.956	0.877	0.913	0.938	0.983	0.960	0.995	0.833	0.891	0.981	0.977	0.978
DeepSeek-R1 (DeepSeek-AI et al., 2025)	0.980	0.730	0.827	0.973	0.990	0.981	0.989	0.821	0.875	0.988	0.980	0.983
GPT-OSS-120B (OpenAI, 2025)	0.984	0.926	0.953	0.956	0.995	0.975	0.976	0.975	0.976	0.992	0.980	0.985
Mistral Large (2407) (AI, 2023)	0.955	0.603	0.724	0.968	0.978	0.972	0.970	0.780	0.848	0.964	0.990	0.976

Table 1: Evaluation of consistent fact sets F' on two datasets (VitaminC (Schuster et al., 2021) / FEVER (Thorne et al., 2018)). Precision (P), recall (R), and F1 are computed with respect to gold consistent subsets. QXR yields cleaner and more complete F' than direct all-at-once LLM judging.

any MUS discovered by Algorithm 1, and I be the number of outer rounds of Algorithm 1 until termination.

Theorem 4.2 (Query Complexity of Algorithm 1). *Algorithm 1 makes at most $I \cdot m \cdot (k \log N)$ oracle calls to LLM O .*

5 Experiments

5.1 Experiment Setting.

Datasets We use VitaminC (Schuster et al., 2021) and FEVER (Thorne et al., 2018) to global consistency by grouping 30 factual statements per example. In VitaminC, each cluster contains 24–28 compatible claims and 2–6 injected contradictions from REFUTES edits, each forming a size-2 ground-truth MUS $\{c^+, c^-\}$. In FEVER, we combine 0–8 REFUTES claims (with evidence) with 14–18 SUPPORTS claims; each refuting claim and its evidence define a ground-truth MUS. Removing either the refuting claim or its evidence yields consistent variants, producing dense contradiction structures that require multi-claim reasoning.

Direct LLM Consistency. We first consider a *direct-LLM baseline* that queries the model once over the entire fact set F , prompting it to return the largest subset $F' \subseteq F$ that is jointly consistent. This corresponds to treating the LLM as an unstructured oracle $O(F)$ that attempts to approximate the ground-truth function $A(F)$ in a single step.

Evaluation. Given an initial fact set F , the model produces a repaired subset F' after removing claims identified as inconsistent. We evaluate the resulting F' against the **gold consistent subset** F_{gold} , defined as the maximal subset of F that contains no injected contradictions (i.e., all ground truth satisfiable claims). Precision, recall, and F1 are computed on the surviving facts:

$$P = \frac{|F' \cap F_{\text{gold}}|}{|F'|}, \quad R = \frac{|F' \cap F_{\text{gold}}|}{|F_{\text{gold}}|}.$$

This directly measures how accurately the model preserves all and only the consistent information, matching the formal objective in Theorem A.1. All experiments use a single oracle call per query ($r = 1$) and the simplest scope setting $C = F$, without majority voting or independence assumptions; such assumptions are used only in the theoretical analysis.

5.2 Results and Analysis.

Table 1 shows consistent gains from MUS-based reasoning. Across all models, QXR yields cleaner consistent subsets F' with substantially higher F1 by avoiding the over-removal seen in direct all-at-once prompting, which often “panic-prunes” large clusters once any conflict is detected. By adaptively isolating minimal conflicts and repairing only what is necessary, QXR preserves nearly all valid information while restoring global consistency. We observe the same trend on a synthetic dataset (Direct: P=0.515, R=0.993, F1=0.644; QXR: P=0.664, R=0.878, F1=0.724), indicating more targeted inconsistency identification (see Appendix E).

All experiments use a zero-shot LLM judge. Evaluating the direct baseline with Chain-of-Thought, decomposition, few-shot, and self-consistency prompting yields the same failure mode: high precision but low recall, showing robustness to prompt design (see Appendix F).

6 Conclusion

We introduced the task of *global fact consistency verification* under noisy LLM oracles, established limits on pairwise sufficiency and query complexity, and proposed an adaptive algorithm that localizes minimal inconsistent subsets and repairs them via hitting-set. On VitaminC clusters, the method improves recall and F1 while preserving high precision, showing that structured querying can turn LLMs into scalable consistency checkers. Future work will extend to larger knowledge graphs and integrate with retrieval and summarization pipelines.

346 Limitations

347 Our theoretical analysis assumes repeated oracle
348 queries can reduce noise under approximate inde-
349 pendence. In practice, however, LLM errors may
350 be systematic rather than random, and repeated
351 queries to the same model do not necessarily im-
352 prove reliability. For this reason, our empirical
353 evaluation does not rely on majority voting or re-
354 peated queries: all experiments use a single LLM
355 call per query ($r = 1$) with the simplest scope set-
356 ting ($C = F$). We treat the independence assump-
357 tion solely as a modeling abstraction for deriving
358 worst-case guarantees.

359 A related limitation is that noise reduction in
360 real systems may require querying *diverse* LLMs
361 rather than repeatedly querying the same model.
362 Exploring ensemble or cross-model consistency
363 checking is a promising direction for future work,
364 but is beyond the scope of this paper.

365 Finally, our experiments focus on moderate-
366 sized fact sets constructed from existing bench-
367 marks. While these settings already expose sub-
368 stantial failures of direct LLM judging, larger and
369 more heterogeneous fact collections—such as those
370 arising in long-context RAG or large-scale knowl-
371 edge extraction—may introduce additional chal-
372 lenges. Designing benchmarks that better capture
373 such regimes remains an open problem.

374 References

- 375 Mistral AI. 2023. Mixtral of experts: release
376 of mixtral 8x7b. [https://mistral.ai/news/
377 mixtral-of-experts](https://mistral.ai/news/mixtral-of-experts).
- 378 Rami Aly and Andreas Vlachos. 2022. Natural logic-
379 guided autoregressive multi-hop document retrieval
380 for fact verification. In *Proceedings of the 2022 Con-
381 ference on Empirical Methods in Natural Language
382 Processing*, pages 6123–6135, Abu Dhabi, United
383 Arab Emirates. Association for Computational Lin-
384 guistics.
- 385 Anthropic. 2025a. Claude 3.7 sonnet and claude
386 code. [https://www.anthropic.com/news/
387 claude-3-7-sonnet](https://www.anthropic.com/news/claude-3-7-sonnet). Accessed: 2025-10-06.
- 388 Anthropic. 2025b. Introducing claude 4. [https://www.
389 anthropic.com/news/claude-4](https://www.anthropic.com/news/claude-4).
- 390 Han Cao, Lingwei Wei, Wei Zhou, and Songlin
391 Hu. 2025. Enhancing multi-hop fact verification
392 with structured knowledge-augmented large language
393 models. *Preprint*, arXiv:2503.08495.
- 394 Ting-Chih Chen, Chia-Wei Tang, and Chris Thomas.
395 2024. *MetaSumPerceiver: Multimodal multi-
396 document evidence summarization for fact-checking*.
397 In *Proceedings of the 62nd Annual Meeting of the*

- Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8742–8757, Bangkok, Thailand. Association for Computational Linguistics. 398
399
400
- Yingjian Chen, Haoran Liu, Yinong Liu, Jinxiang Xie, Rui Yang, Han Yuan, Yanran Fu, Peng Yuan Zhou, Qingyu Chen, James Caverlee, and Irene Li. 2025. *Graphcheck: Breaking long-term text barriers with extracted knowledge graph-powered fact-checking*. *Preprint*, arXiv:2502.16514. 401
402
403
404
405
406
- Marie-Catherine de Marneffe, Anna N. Rafferty, and Christopher D. Manning. 2008. *Finding contradictions in text*. In *Proceedings of ACL-08: HLT*, pages 1039–1047, Columbus, Ohio. Association for Computational Linguistics. 407
408
409
410
411
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. *Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning*. *Preprint*, arXiv:2501.12948. 412
413
414
415
416
417
418
419
- Bishwamitra Ghosh, Sarah Hasan, Naheed Anjum Arafat, and Arijit Khan. 2025. *Logical consistency of large language models in fact-checking*. In *The Thirteenth International Conference on Learning Representations*. 420
421
422
423
424
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. 2025. *A survey on llm-as-a-judge*. *Preprint*, arXiv:2411.15594. 425
426
427
428
429
430
- Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. *A survey on automated fact-checking*. *Transactions of the Association for Computational Linguistics*, 10:178–206. 431
432
433
434
- Zhaochen Hong, Haofei Yu, and Jiaxuan You. 2025. *ConsistencyChecker: Tree-based evaluation of LLM generalization capabilities*. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 33039–33075, Vienna, Austria. Association for Computational Linguistics. 435
436
437
438
439
440
441
- Qisheng Hu, Quanyu Long, and Wenya Wang. 2025. *Decomposition dilemmas: Does claim decomposition boost or burden fact-checking performance?* In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6313–6336, Albuquerque, New Mexico. Association for Computational Linguistics. 442
443
444
445
446
447
448
449
450
- Ulrich Junker. 2004. *Quickxplain: Preferred explanations and relaxations for over-constrained problems*. pages 167–172. 451
452
453

454	Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9332–9346, Online. Association for Computational Linguistics.	511
455		512
456		513
457		514
458		
459		515
460		516
461	Sujit Kumar, Rohan Jaiswal, Mohit Ram Sharma, and Sanasam Ranbir Singh. 2023. Multiset dual summarization for incongruent news article detection . In <i>Proceedings of the 20th International Conference on Natural Language Processing (ICON)</i> , pages 779–790, Goa University, Goa, India. NLP Association of India (NLP AI).	517
462		518
463		519
464		520
465		
466		521
467		522
468	Jimmy Lee and K. Leung. 2010. A stronger consistency for soft global constraints in weighted constraint satisfaction . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 24:121–127.	523
469		524
470		525
471		526
472	Patrick 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. 2021. Retrieval-augmented generation for knowledge-intensive nlp tasks . <i>Preprint</i> , arXiv:2005.11401.	
473		527
474		528
475		529
476		530
477		531
478	Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024. Llms-as-judges: A comprehensive survey on llm-based evaluation methods . <i>Preprint</i> , arXiv:2412.05579.	532
479		533
480		534
481		535
482		536
483	Tianyu Liu, Jirui Qi, Paul He, Arianna Bisazza, Mrinmaya Sachan, and Ryan Cotterell. 2025. Pointwise mutual information as a performance gauge for retrieval-augmented generation . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1628–1647, Albuquerque, New Mexico. Association for Computational Linguistics.	537
484		538
485		539
486		540
487		541
488		
489		542
490		543
491		
492	OpenAI. 2025. Introducing gpt-oss . https://openai.com/index/introducing-gpt-oss/ .	544
493		545
494	Liangming Pan, Xiaobao Wu, Xinyuan Lu, Anh Tuan Luu, William Yang Wang, Min-Yen Kan, and Preslav Nakov. 2023. Fact-checking complex claims with program-guided reasoning . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6981–7004, Toronto, Canada. Association for Computational Linguistics.	546
495		547
496		548
497		549
498		550
499		551
500		552
501		
502	Subhey Sadi Rahman, Md. Adnanul Islam, Md. Mahub Alam, Musarrat Zeba, Md. Abdur Rahman, Sadia Sultana Chowa, Mohaimenul Azam Khan Raiaan, and Sami Azam. 2025. Hallucination to truth: A review of fact-checking and factuality evaluation in large language models . <i>Preprint</i> , arXiv:2508.03860.	553
503		554
504		555
505		556
506		557
507		558
508	Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. Get your vitamin C! robust fact verification with contrastive evidence . In <i>Proceedings of the 2021</i>	
509		559
510		560
		561
		562
		563
		564
		565
		511
		512
		513
		514
		515
		516
		517
		518
		519
		520
		521
		522
		523
		524
		525
		526
		527
		528
		529
		530
		531
		532
		533
		534
		535
		536
		537
		538
		539
		540
		541
		542
		543
		544
		545
		546
		547
		548
		549
		550
		551
		552
		553
		554
		555
		556
		557
		558
		559
		560
		561
		562
		563
		564
		565

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging LLM-as-a-judge with MT-bench and chatbot arena](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Lianghui Zhu, Xinggang Wang, and Xinlong Wang. 2025. [Judgelm: Fine-tuned large language models are scalable judges](#). *Preprint*, arXiv:2310.17631.

A Theorems and Proofs

Let $F = \{f_1, \dots, f_N\}$ be a finite fact set, $C = \{C_1, \dots, C_m\}$ with $C_i \subseteq F$ be a collection of constraint scopes, and a perfect oracle $O : 2^F \rightarrow \{\text{cons}, \text{incons}\}$ that returns whether a subset of facts is jointly consistent. Throughout the proofs, we will use two standard properties.

- **Monotonicity:** If $U \subseteq S \subseteq F$ and $O(U) = \text{incons}$ then $O(S) = \text{incons}$.
- **Existence of MUS:** If $O(S) = \text{incons}$ and S is finite then S contains a MUS $U \subseteq S$.

These follow from finiteness and the definition of minimality.

A.1 Objective Equivalence and Reduction to Hitting Set

Theorem A.1. *Objective Equivalence Maximizing coverage*

$$\max_{F' \subseteq F} |F'| \quad \text{s.t.} \quad A(F' \cap C_i) = \text{cons.} \quad \forall i \in [m] \quad (3)$$

is equivalent to minimizing deletions

$$\min_{R \subseteq F} |R| \quad \text{s.t.} \quad A((F \setminus R) \cap C_i) = \text{cons.}, \quad \forall i \in [m] \quad (4)$$

Proof. Define a bijection between solutions $R = F \setminus F'$ and $F' = F \setminus R$. Then

$$|F'| = |F| - |R|. \quad (5)$$

Hence, maximizing $|F'|$ is equivalent to minimizing $|R|$. Under this bijection, the feasibility constraints are identical. \square

Let $\mathcal{A}_{\text{incons}}$ be the family of all MUSes possible in C . Formally, based on Definition 4.1, we define

$$\mathcal{A}_{\text{incons}} = \{U \in C_j : j \in [m], U \text{ is a MUS w.r.t. } A\}.$$

Theorem A.2. *A set $R \subseteq F$ is feasible iff it is a hitting set for $\mathcal{A}_{\text{incons}}$.*

Proof. In one direction, assume $R \subseteq F$ is feasible for minimum deletion. Suppose for contradiction, there exists $U \in \mathcal{A}_{\text{incons}}$ with $R \cap U = \emptyset$. By definition of $\mathcal{A}_{\text{incons}}$, there is some i with $U \subseteq C_i$, and thus $U \subseteq (F \setminus R) \cap C_i$. Since $A(U) = \text{incons}$ by monotonicity $(F \setminus R) \cap C_i$ would be inconsistent, contradiction to Equation (4). Hence, $R \cap U \neq \emptyset$ for all $U \in \mathcal{A}_{\text{incons}}$; i.e., R is a hitting set.

For the other direction, assume R is a hitting set for $\mathcal{A}_{\text{incons}}$, i.e. $R \cap U \neq \emptyset$ for all $U \in \mathcal{A}_{\text{incons}}$. Suppose for contradiction that R is not feasible; then there exists some i with $A((F \setminus R) \cap C_i) = \text{incons}$. By the existence-of-MUS property, $(F \setminus R) \cap C_i$ contains a MUS $U \subseteq (F \setminus R) \cap C_i$. Then $U \in \mathcal{A}_{\text{incons}}$ but $U \cap R = \emptyset$, contradicting that R hits all MUSes. Hence $A((F \setminus R) \cap C_i) = \text{cons}$ for all i , i.e. R is feasible. \square

A.2 Pairwise Checks Are Insufficient for Global Consistency

Theorem A.3 (Pairwise Insufficiency). *Even with access to A , there exist F of size $N \geq 3$ such that all pairs are consistent, but F is inconsistent globally.*

Proof. Given A we now want to show pairwise consistency does not imply global consistency. Let the universe be boolean assignments to variables $A, B, C \in \{0, 1\}$. Consider the three facts $f_1 : A \oplus B = 1$, $f_2 : B \oplus C = 1$, $f_3 : C \oplus A = 1$. Then, any two pairs can be satisfied, for example f_1, f_2 are satisfied with $A = 1, B = 0, C = 1$. However, this is jointly unsatisfiable, from f_1 and f_2 : $A = \neg B$ and $C = \neg B$ hence $A = C$, then $C \oplus A = 0$ which contradicts f_3 . Hence, $f_1, \wedge f_2, \wedge f_3$ is inconsistent. We can think of this as a graph colouring problem, for example the constraints $A \neq B, B \neq C, C \neq A$ requires a 2-colouring of a 3-cycle, which is impossible. \square

A.3 Soundness of Algorithm 1

Theorem A.4 (Soundness under Perfect Oracle). *Assume that the LLM oracle O has zero error ($\alpha = 0, \beta = 0$), then for every $i \in [m]$, the retained subset $F' \triangleq F \setminus R$ satisfies*

$$A(F' \cap C_i) = \text{cons.} \quad (6)$$

Proof. If O is perfect, then it should match the ground-truth A on all inputs. Suppose for contradiction that $O(F' \cap C_i) = \text{incons}$ for some $i \in [m]$. By MUS existence, there is a MUS $U \subseteq F' \cap C_i$. Then, $U \subseteq C_i$, so $U \in \mathcal{A}_{\text{incons}}$. We know $U \subseteq F \setminus R$ so we know $U \cap R = \emptyset$ contradicting R hits every $U \in \mathcal{A}_{\text{incons}}$. Therefore, $O(F' \cap C_i) = \text{cons}$ for all i . \square

In practical scenarios, the MUSes $\mathcal{O}_{\text{incons}}$ extracted by Algorithm 1 may contain errors, as the procedure depends on a noisy oracle O . However, F' generated by in Algorithm 1 is still sound with respect to the MUSes found.

Theorem A.5 (Practical Soundness w.r.t. Extracted Conflicts). *Let $\mathcal{O}_{\text{incons}}$ be the set of MUSes actually extracted by the Algorithm 1 under O . Let $R \subseteq F$ satisfy $R \cap U \neq \emptyset$ for all $U \in \hat{\mathcal{U}}$. The retained set $F \setminus R$ is consistent with respect to the extracted conflicts $\mathcal{O}_{\text{incons}}$.*

Proof. For every $U \in \mathcal{O}_{\text{incons}}$, $R \cap U \neq \emptyset$ implies $U \not\subseteq F \setminus R$. Thus none of the discovered MUSes in $\mathcal{O}_{\text{incons}}$ remain after removal. \square

A.4 Error Reduction under Repetition

Let the true perfect oracle answer for a query be $Y \in \{\text{cons}, \text{incons}\}$. A noisy oracle O is repeatedly queried r times on the same set S . Each repetition returns $\hat{Y}_t \in \{\text{cons}, \text{incons}\}$, $t = 1, \dots, r$. For each repetition and conditioning on the true label Y , we get

$$\Pr(\hat{Y}_t \neq Y | Y) \leq \varepsilon \quad \varepsilon \triangleq \max\{\alpha, \beta\} < \frac{1}{2} \quad (7)$$

We assume independence across repetitions, e.g., $\{\hat{Y}_t\}_{t=1}^r$ are independent conditioned on Y . We also assume r is odd.

Theorem A.6 (Error Reduction Under Repetition). *Let \mathcal{O} be an (α, β) -noisy subset-consistency oracle with $\max\{\alpha, \beta\} < \frac{1}{2}$. If each query is evaluated r times independently and aggregated by majority vote, the effective error rate per aggregated call is at most $\exp(-2r\gamma^2)$ where $\gamma = \frac{1}{2} - \max\{\alpha, \beta\}$.*

Proof. The majority vote is wrong iff at least half of the repetitions are wrong, using a conservative threshold we get

$$\{\hat{Y}^{\text{maj}} \neq Y\} \subseteq \{S_r \geq r/2\}$$

where $X_t = \mathbb{I}\{\hat{Y}_t \neq Y\}$ and $S_r = \sum_{t=1}^r X_t$. So we can apply Hoeffding's inequality to S_r (sum of independent Bernoulli variables with means $\leq \varepsilon \triangleq \max\{\alpha, \beta\} < \frac{1}{2}$). For any $a > 0$, we have

$$\Pr(S_r - \mathbb{E}[S_r] \geq a) \leq \exp\left(-\frac{2a^2}{r}\right)$$

here $\mathbb{E}[S_r] = \sum_{t=1}^r \mathbb{E}[X_t] \leq r\varepsilon$. Set $a = r(\frac{1}{2} - \varepsilon) = r\gamma$. Then,

$$\begin{aligned} \Pr(S_r \geq r/2) &\leq \Pr(S_r - \mathbb{E}[S_r] \geq r(\frac{1}{2} - \varepsilon)) \\ &\leq \exp(-2r\gamma^2). \end{aligned}$$

Algorithm 2: QuickXplain (QX) for MUS Extraction

Input: Oracle O ;

Candidate set S ;

Background B (assumed consistent)

Output: Subset-minimal inconsistent set

$$\Delta \subseteq S \text{ (or } \emptyset)$$

```

1 if  $S = \emptyset$  then
2   return  $\emptyset$ 
3 if  $O(B \cup S) = \text{cons}$  then
4   return  $\emptyset$ 
5 if  $|S| = 1$  then
6   return  $S$ 
7 Split  $S$  into two (nearly) equal parts  $S_1, S_2$ ;
8  $\Delta_1 \leftarrow \text{QX}(O, S_1, B \cup S_2)$ ;
9  $\Delta_2 \leftarrow \text{QX}(O, S_2, B \cup \Delta_1)$ ;
10 return  $\Delta_1 \cup \Delta_2$ 

```

Combining together we get per-call bound

$$\Pr(\hat{Y}^{\text{maj}} \neq Y) \leq \exp(-2r\gamma^2).$$

\square

B Hitting Set Approximation

Lemma B.1 (Greedy Hitting Set Approximation (Vazirani, 2001)). *Let \mathcal{U} be a family of m MUSes (conflict sets) over facts F . The greedy hitting set algorithm returns a hitting set H satisfying $|H| \leq (1 + \ln m) |H^*|$, where H^* is an optimal (minimum-cardinality) hitting set of \mathcal{U} .*

C QuickXplain

QuickXplain (QX) is a classic divide-an-conquer method for localizing a minimal unsatisfiable subset (MUS) from an inconsistent set S under a consistency oracle $O : 2^F \rightarrow \{\text{cons}, \text{incons}\}$ (Junker, 2004). The algorithm adaptively queries subsets to find a subset-minimal inconsistent core with logarithmic depth in $|S|$. We provide a pseudocode in Algorithm 2. QX narrows the inconsistent subset by recursively testing halves of S . If $B \cup S_1$ is inconsistent, the conflict lies in S_1 ; otherwise it lies in S_2 . Recursion stops once singleton conflicts are reached, yielding a subset-minimal inconsistent set. For a perfect oracle, QX requires $\mathcal{O}(k \log |S|)$ queries to locate a conflict of size k .

D Pairwise/NLI Baselines are not comparable at scale

Methods that decide consistency by evaluating all sentence pairs require $\binom{N}{2}$ oracle calls. For typi-

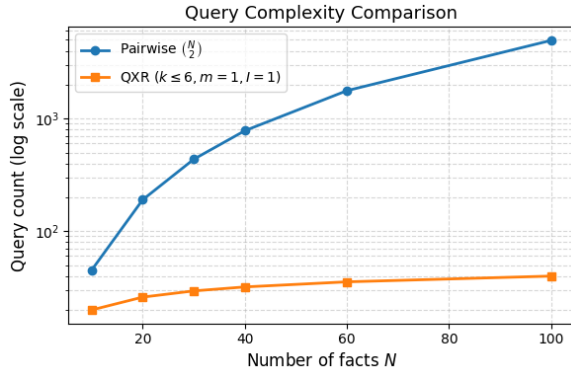


Figure 1: Scaling of query counts with number of facts N . Pairwise checking is quadratic, while QXR scales polylogarithmically

cal sizes of facts e.g., $N \in [30, 100]$, that implies 435 – 4950 calls *per instance*. Under noisy LLM oracles, each decision further requires r repetitions (majority vote). Our QXR algorithm performs at most $I \cdot m \cdot (k \log N)$ number of queries where k is the MUS maximum size (empirically small), m is the constraint scope (if no constraints or in general $m = 1$) and I is the number of outer rounds. In our experiments $k \leq 6$, $m = 1$, we show a plot how complexity scales in Figure 1. We also illustrate common baselines and why they do not scale in Table 2.

E Synthetic Dataset

Entities and predicates. We sample entity names for PERSON, ORG, LOC, EVENT, and ANIMAL. Facts are rendered from a small predicate bank with negation support, including unary categories (IsTiger, IsDog, IsActor, IsPolitician, IsAnimal), binary relations (WorksFor, LocatedIn), and temporal precedence (Before). Numeric paraphrases (AtLeastCases, AtMostCases) provide clutter.

Planted MUS patterns. We inject one or more MUSes per instance:

- size-2 contradiction ($A, \neg A$),
- size-3 temporal cycle,
- size-3 “exactly one” parity conflict consisting of a rule sentence (*Exactly one of X, Y, Z holds*) plus two of the three facts.

Distractors. We add on-topic true facts and a small fraction of off-topic false facts (e.g., *does not work for, is not located in*) to create realistic clutter.

Gold Annotations. For each instance we store the gold MUS family $\mathcal{U}_{\text{gold}}$ as lists of sentence IDs. The gold consistent subset F_{gold} is defined as the maximal consistent subset obtained by removing a minimum hitting set over $\mathcal{U}_{\text{gold}}$ (greedy approximation).

F Prompting Analysis

All main experiments in the paper use a zero-shot LLM-judge setting. To assess sensitivity to prompt design, we additionally evaluated the direct-judge baseline on VitaminC using several advanced prompting strategies. Across all prompt styles, we observe the same qualitative trend: high precision but substantially lower recall due to over-removal of consistent facts. These experiments were run on a subset of the dataset; for reference, we also restate the original zero-shot baseline and QXR results from the full evaluation.

Claude-3.7 (Sonnet)

Prompting Style	P	R	F1
Zero-shot	0.986	0.681	0.797
Chain-of-Thought	0.973	0.684	0.792
Decomposition	0.987	0.626	0.755
Few-shot	0.975	0.658	0.775
Self-consistency	0.986	0.547	0.694
Original zero-shot baseline	0.979	0.854	0.909
Original QXR	0.956	0.975	0.965

Claude-4 (Sonnet)

Prompting Style	P	R	F1
Zero-shot	0.989	0.928	0.954
Chain-of-Thought	0.987	0.908	0.942
Decomposition	0.677	0.608	0.638
Few-shot	0.992	0.891	0.932
Self-consistency	0.992	0.768	0.860
Original zero-shot baseline	0.956	0.877	0.913
Original QXR	0.938	0.983	0.960

DeepSeek-R1

Prompting Style	P	R	F1
Zero-shot	0.993	0.860	0.910
Chain-of-Thought	0.996	0.850	0.907
Decomposition	0.996	0.882	0.930
Few-shot	0.882	0.761	0.813
Self-consistency	0.987	0.718	0.817
Original zero-shot baseline	0.980	0.730	0.827
Original QXR	0.973	0.990	0.981

GPT-OSS-120B

Prompting Style	P	R	F1
Zero-shot	0.994	0.876	0.920
Chain-of-Thought	0.995	0.851	0.903
Decomposition	0.954	0.816	0.865
Few-shot	0.997	0.896	0.932
Self-consistency	0.996	0.475	0.624
Original zero-shot baseline	0.984	0.926	0.953
Original QXR	0.956	0.995	0.975

Method	Description	Complexity
Pairwise NLI (FEVER-style)	Run NLI on every pair of claims, removing any node involved in a contradiction edge.	$O(N^2)$
Transitive Closure / Entailment Graph	Build a full entailment-contradiction graph and perform reasoning (e.g., SAT solving or closure).	$> O(N^2)$
LLM-as-NLI	Prompt an LLM with two sentences (e.g., "Does A contradict B?").	$O(N^2)$
Multi-premise NLI	Treat the entire set of claims as premises and ask if they jointly entail a hypothesis.	$O(1)$
Clustered Pairwise	For each claim, compare only to its top- K nearest neighbors (embedding-based).	$O(NK)$

Table 2: Common NLI-style baselines for consistency checking and their computational complexity. Pairwise and entailment-graph methods grow quadratically with the number of claims, making them infeasible for large clusters. Multi-premise NLI corresponds to our baseline.

Mixtral-8×7B

Prompting Style	P	R	F1
Zero-shot	0.947	0.568	0.701
Chain-of-Thought	0.997	0.570	0.700
Decomposition	0.991	0.535	0.673
Few-shot	0.238	0.142	0.175
Self-consistency	0.990	0.522	0.674
Original zero-shot baseline	0.955	0.603	0.724
Original QXR	0.968	0.978	0.972

{facts_block}

CRITICAL: Return ONLY a Python list using the FULL TEXT of each fact.
<answer>["fact1", "fact2", ...]</answer>

Direct baseline (Chain-of-Thought prompting).

We encourage structured reasoning while constraining the output format:

Given the following facts, some may contradict. Find all mutually consistent facts.

Facts:
{facts_block}

Think step-by-step:
1. Identify pairs of facts that contradict each other.
2. For each contradiction, decide which fact to keep.
3. Return the consistent subset.

CRITICAL: In the <answer> tag, return the FULL TEXT of each fact, NOT numbers.
Example: <answer>["The sky is blue", "Grass is green"]</answer>

Direct baseline (Decomposition prompting).

The task is decomposed into detection and resolution:

Task: Select a mutually consistent subset of the following facts.

Facts:
{facts_block}

Step 1 - Identify contradicting facts.
Step 2 - Decide which facts to keep.
Step 3 - Output the consistent subset.

CRITICAL: Return ONLY a Python list with the FULL TEXT of each fact.
<answer>["full fact text 1", "full fact text 2", ...]</answer>

G Prompts

This appendix lists the exact prompts used in our experiments. All prompts were intentionally kept simple and symmetric across methods to isolate algorithmic effects rather than prompt engineering. For all baselines that output a subset of facts, models are required to return the *full text* of each retained fact (not indices) as a Python list enclosed in an <answer> tag.

Subset-consistency oracle (QXR). Given a queried subset of facts (optionally with a background set B), we ask the LLM to judge whether all statements can be true simultaneously:

Factual statements. Some may contradict.

{bg}Statements:
{facts_block}

Respond ONLY:
- CONSISTENT
- INCONSISTENT

Answer:

Direct baseline (zero-shot). The LLM is asked to return a mutually consistent subset of facts:

Given the following factual statements, some may contradict.

Return a Python list of facts that are mutually consistent, meaning all returned facts can be true at the same time.

Facts:

869 **Direct baseline (Few-shot prompting).** We pro-
870 vide two illustrative examples followed by the tar-
871 get instance:

872 Given facts that may contradict, return
873 the subset that is mutually consistent.

874 Example 1:

875 - The sky is blue
876 - The sky is red
877 - Grass is green
878 <answer>["The sky is blue", "Grass is
879 green"]</answer>

881 Example 2:

882 - Paris is in France
883 - Paris has 2 million people
884 - Paris has 10 million people
885 <answer>["Paris is in France", "Paris
886 has 2 million people"]</answer>

887 Now solve:

888 Facts:
889 {facts_block}

890 <answer>

894 **Direct baseline (Self-consistency prompting).**

895 We sample multiple outputs using the same prompt
896 and aggregate by majority vote over selected facts:

897 Given the following facts, some may
898 contradict.
899 Return ONLY a Python list of facts that
900 are mutually consistent.

901 Facts:
902 {facts_block}

903 <answer>["fact1", "fact2",
904 ...]</answer>