

# A Benchmark for the Detection of Metalinguistic Disagreements between LLMs and Knowledge Graphs

Bradley P. Allen<sup>1,\*</sup>, Paul T. Groth<sup>1</sup>

<sup>1</sup>University of Amsterdam, Science Park 900, 1098 XH Amsterdam, The Netherlands

## Abstract

Evaluating large language models (LLMs) for tasks like fact extraction in support of knowledge graph construction frequently involves computing accuracy metrics using a ground truth benchmark based on a knowledge graph (KG). These evaluations assume that errors represent factual disagreements. However, human discourse frequently features *metalinguistic disagreement*, where agents differ not on facts but on the meaning of the language used to express them. Given the complexity of natural language processing and generation using LLMs, we ask: do metalinguistic disagreements occur between LLMs and KGs? Based on an investigation using the T-REx knowledge alignment dataset, we hypothesize that metalinguistic disagreement does in fact occur between LLMs and KGs, with potential relevance for the practice of knowledge graph engineering. We propose a benchmark for evaluating the detection of factual and metalinguistic disagreements between LLMs and KGs. An initial proof of concept of such a benchmark is available on Github.

## Keywords

large language models, knowledge graphs, fact checking, metalinguistic disagreement

## 1. Introduction

Recent years have seen a surge of interest in the use of LLMs for purposes of knowledge engineering [1]. LLMs are being used to perform text classification, sentiment analysis, and natural language inference, exploiting next-token prediction to generate text that can be transformed into the type of symbolic outputs normally produced in these tasks [2]. Increasing emphasis is being placed on the use of LLMs in knowledge graph construction [3]. The results have been encouraging, but a major concern that has emerged is the impact of hallucination, which is defined as the presence of factually incorrect or unjustified assertions in the output of LLMs [4, 5]. Benchmarks such as SHROOM [6] and WildHallucinations [7] have been developed to evaluate the ability to detect hallucination when it occurs in LLM output.

A number of mechanisms have been proposed to mitigate hallucination in LLMs through the use of knowledge from a variety of sources, including natural language text, KGs, and rules, to ground [8] an LLM. Retrieval-augmented generation (RAG) is a specific version of this approach that has attracted a great deal of interest, particularly in the context of commercial applications [9]. Such *knowledge-enhanced LLMs* [10] show improvements in the performance of natural language understanding and generation tasks. However, even with such improvements, knowledge-enhanced LLMs still produce errors as measured using common evaluation metrics (e.g. F1 measures for classification). These evaluation metrics are calculated by measuring the difference between an LLM's output and ground truth as provided in fact checking benchmarks such as LAMA [11], KAMEL [12], and FActScore [13].

The errors reported by these metrics are typically assumed to stem from disagreement about facts. But there is another way in which these differences can arise. *Metalinguistic disagreement* [14, 15, 16] occurs when people argue about the meaning or use of words rather than about facts or ideas. In contrast, a factual disagreement is about what is actually true in the world. Examples of factual disagreement are debating whether a tomato is healthier than an apple, or debating whether Sarah is taller than John; in

*The 23rd International Semantic Web Conference, November 11–15, 2025, Baltimore, MD*

\*Corresponding author.

✉ b.p.allen@uva.nl (B. P. Allen); p.t.groth@uva.nl (P. T. Groth)

🌐 <https://www.bradleypallen.org/> (B. P. Allen); <https://pgroth.com/> (P. T. Groth)

🆔 0000-0003-0216-3930 (B. P. Allen); 0000-0003-0183-6910 (P. T. Groth)



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

contrast, examples of metalinguistic disagreement are arguing whether a tomato should be called a fruit or a vegetable, or arguing about what height qualifies as “tall” when describing a person.

Consider the following scenario: a knowledge-enhanced LLM generates an output that contradicts ground truth provided by a KG. This is used as evidence that the LLM has committed a factual error in its output. However, in producing its output, the knowledge-enhanced LLM has provided a rationale that indicates that there is a disagreement about the meaning of a term that has led to the output. Can this occur in practice? Our hypothesis is that it does.

Why would this matter? Factual disagreements can be resolved through knowledge graph refinement [17] or through few-shot in-context learning that provides the correct facts to the LLM; however, metalinguistic disagreements may require ontology engineering to address representational issues with a knowledge graph, or the engineering of prompts that incorporate intensional definitions in natural language of concepts for an LLM [18]. Data governance [19] also acknowledges the importance of establishing metalinguistic agreement of intensional definitions of concepts and relations in natural language and their realization in databases and database schemas; for example, the FAIR principles [20, 21] specifically urge clear documentation of metadata aligning natural language concepts and metadata in scientific data resources. We therefore argue that distinguishing factual from metalinguistic disagreement between LLMs and KGs is relevant to the practice of knowledge graph and ontology engineering.

## 2. Evidence for the occurrence of metalinguistic disagreement in LLMs

To test our hypothesis that metalinguistic disagreement is a detectable phenomenon, we conducted a simple experiment by fact checking a set of knowledge graph triples aligned with natural language text using an LLM, and then estimating the rate at which metalinguistic disagreement occurs when the LLM determines the triple is not true.

We randomly sampled 100 Wikipedia abstracts from the 10,000 document sample provided in the T-REx dataset, a dataset of large scale alignments between Wikipedia abstracts and Wikidata triples. T-REx has been widely used in the evaluation of LLM-based fact checking and extraction for knowledge graphs [22]. From the total set of triples aligned with the documents in that sample, we then sampled 250 triples. We then defined a zero-shot chain-of-thought classifier [23] to assign a truth value to an aligned triple, providing the Wikidata abstract with which it is aligned as context in the LLM prompt [24]. The classifier was executed to obtain a rationale and a truth value for each of the 250 sampled triples and aligned abstracts, and each result was then processed by a second zero-shot chain-of-thought classifier (using gpt-4o-2024-05-13) acting as an LLM-as-a-judge [25], to classify whether the truth-value-assigning classifier’s rationale indicated a metalinguistic disagreement. Processing required a total of 2 inference API calls per alignment, per LLM. Evaluations whose statistics are reported below were conducted during the period from 1 July 2024 to 8 July 2024. Costs incurred through calls to language model APIs totalled less than \$100 USD. Code and data used in the experiments are available in a Github repository<sup>1</sup>.

As shown in Table 1, over the 9 LLMs evaluated, false negative rates over the 250 sampled T-REx triples ranged between 0.104 and 0.504 with a mean of 0.246, and the rate of metalinguistic disagreements between the classifier and Wikidata (i.e., the number of detected metalinguistic disagreements divided by the number of evaluated alignments) ranged between 0.04 and 0.264 with a mean of 0.097. To illustrate the nature of disagreements detected, Table 2 shows two examples of false negatives from the experiment which exhibit metalinguistic disagreement. The complete set of rationales and classifications is available in the Github repository.

---

<sup>1</sup><https://github.com/bradleypallen/trex-metalinguistic-disagreement>

LLM	FN	FNR	metalinguistic disagreements	MDR
gpt-4o-2024-05-13	26	0.104	10	0.040
gpt-4-0125-preview	33	0.132	16	0.064
claude-3-haiku-20240307	42	0.168	11	0.044
claude-3-opus-20240229	51	0.204	14	0.056
claude-3-5-sonnet-20240620	52	0.208	16	0.064
meta-llama/Meta-Llama-3-70b-Instruct	71	0.284	28	0.112
gpt-3.5-turbo	75	0.300	29	0.116
mistralai/Mixtral-8x7B-Instruct-v0.1	78	0.312	29	0.116
mistralai/Mistral-7B-Instruct-v0.3	126	0.504	66	0.264

**Table 1**

Performance by LLM reporting the number of false negatives (FN); the false negative rate (FNR); and the number and rate of metalinguistic disagreements (MDR) detected by the LLM-as-a-judge.

Wikidata triple	rationale for negative finding
James II of England (Q126188) followed by (P156) Mary II of England (Q130812)	Based on the information provided in the hypothetical world, the knowledge graph triple “Subject: James II of England, Predicate: followed by, Object: Mary II of England” is false. The passage states that after James II was deposed in 1688, he was replaced by his daughter Mary II, ruling jointly with her husband William III. This means that Mary II did not immediately follow James II, but rather replaced him as the monarch. The predicate “followed by” implies that Mary II succeeded James II directly, which is not the case.
chocolate (Q195) made from material (P186) sugar (Q11002)	Based on the information provided, the knowledge graph triple “Subject: chocolate, Predicate: made from material, Object: sugar” is false. The passage clearly states that chocolate is made from the seeds of the cacao tree, specifically Theobroma cacao, and that it is roasted and ground to produce the chocolate. There is no mention of sugar being a material that chocolate is made from. In fact, the passage states that sugar is added to chocolate to create sweet chocolate, but it is not a primary ingredient that chocolate is made from. Therefore, the given knowledge graph triple is false.

**Table 2**

Two examples of rationales (generated using claude-3-haiku-20240307) exhibiting metalinguistic disagreement. In both cases, the disagreement is around the meaning of the predicate. For the triple with James II of England as the subject, the LLM states that the historical facts presented in the context are better represented using “replaced by” (P1366) than “followed by” (P156); for the triple with chocolate as a subject, the LLM states that “made from material” (P186) is not an appropriate relation when the object is not a primary ingredient.

### 3. Proposed benchmark

We argue that the above results suggest that that metalinguistic disagreement between knowledge graphs and LLMs can occur during fact-checking tasks. However, there are some significant shortcomings in the above approach:

- **Lack of human validation.** The detection of metalinguistic disagreement relies on using an LLM-as-a-judge, which may not be a reliable substitute for human judgment [26, 27]. This introduces the possibility that the detected “disagreements” are artifacts of how different LLMs process and generate language, rather than true metalinguistic disagreements. Human review at scale is needed to validate the results. Without this, it’s difficult to determine if what the LLMs identify as metalinguistic disagreements align with human judgments.
- **Possible conflation with other error types.** What’s interpreted as metalinguistic disagreement could potentially be other types of errors or inconsistencies in LLM outputs, such as hallucinations or context misinterpretations.
- **Limited sample size.** The experiment uses a relatively small sample of 250 triples. A larger-scale study is needed to draw more robust conclusions.

We argue that by creating a benchmark metalinguistic disagreement detection dataset that addresses

these limitations, we could more confidently assess the occurrence and nature of metalinguistic disagreements in LLM-based fact-checking. This would provide a stronger foundation for investigating our hypothesis and advancing our understanding of how LLMs interpret and disagree about meaning in knowledge graph engineering contexts.

Specific requirements for such a benchmark include:

- **Human-annotated examples.** A set of fact-checking instances annotated by human experts to identify clear cases of metalinguistic disagreement, factual disagreement, and agreement. This would serve as a gold standard for evaluation.
- **Inter-annotator agreement metrics.** Support the evaluation of system performance using inter-annotator agreement metrics that incorporate knowledge graph ground truth and human annotations to measure the degrees of inter-agent factual and metalinguistic agreement.
- **Multiple knowledge graph sources.** Use triples from different knowledge graphs spanning multiple knowledge domains to account for variations in how relations and concepts are defined across sources, and to test if metalinguistic disagreements are more prevalent in certain areas.
- **Contextual information.** Provide relevant context for each fact-checking instance, similar to the Wikipedia abstracts used by T-REx.
- **Examples with ambiguity, temporal aspects, and gradable predicates.** Deliberately include examples with potential for ambiguity or multiple interpretations to probe the boundaries of metalinguistic disagreement, examples where the truth value of a statement might change over time, to explore how temporal context affects metalinguistic understanding, and examples with gradable predicates (e.g., "tall," "fast") that might be more prone to metalinguistic disagreement.
- **Negative examples.** Include clear cases where no metalinguistic disagreement should occur, to test for false positives.

As an initial next step towards this objective, we plan to extend the dataset used in the initial experiments described above in a manner similar to that used in the design and implementation of the SHROOM hallucination detection benchmark [6, 28], through crowdsourcing to incorporate human annotation and increasing the size of the sample of knowledge alignments from the T-REx dataset. Human annotators will be presented with a summary of a Wikipedia page and a statement generated from the Wikidata knowledge graph triple for each alignment, and the annotator must indicate if they disagree with the statement, and if so, whether they disagree on the factuality of the statement or the meaning of any of the terms used in the statement.

In conclusion, we anticipate that such a benchmark can not only shed light on the nature and frequency of metalinguistic disagreements between LLMs and KGs, but also contribute to the ongoing debate about LLMs' capacity for generating meaningful statements. Some have argued that LLMs are incapable of understanding meaning in the way humans do [29]. Others are exploring ways in which LLMs might be capable of at least some limited or partial forms of meaning as a consequence of either the model's pre-training or its grounding through in-context learning [30, 31, 32, 33, 34]. We believe that the proposed benchmark can contribute to a more nuanced view of the epistemic status of LLMs relative to KGs based on two-component semantics [35, 36, 37], and support experimental work in determining whether or not LLMs can generate meaningful statements or be claimed to have beliefs [38, 39].

## Acknowledgements

This work was partially supported by EU's Horizon Europe research and innovation programme within the ENEXA project (grant Agreement no. 101070305). The authors wish to thank Frank van Harmelen, Levin Hornischer, Filip Ilievski, Jan-Christoph Kalo, Aybüke Özgün, Lise Stork, and Klim Zaporjets for discussions and suggestions that have been invaluable in refining this work.

## References

- [1] B. P. Allen, L. Stork, P. Groth, Knowledge Engineering Using Large Language Models, *Transactions on Graph Data and Knowledge* 1 (2023) 3:1–3:19. URL: <https://drops.dagstuhl.de/entities/document/10.4230/TGDK.1.1.3>. doi:10.4230/TGDK.1.1.3.
- [2] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, G. Neubig, Pre-Train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing, *ACM Comput. Surv.* 55 (2023). URL: <https://doi.org/10.1145/3560815>. doi:10.1145/3560815.
- [3] E. Koutsiana, J. Walker, M. Nwachukwu, A. Meroño-Peñuela, E. Simperl, Knowledge Prompting: How Knowledge Engineers Use Large Language Models, *arXiv preprint arXiv:2408.08878* (2024).
- [4] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, P. Fung, Survey of hallucination in natural language generation, *ACM Computing Surveys* 55 (2023) 1–38.
- [5] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, et al., A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions, *arXiv preprint arXiv:2311.05232* (2023).
- [6] T. Mickus, E. Zosa, R. Vázquez, T. Vahtola, J. Tiedemann, V. Segonne, A. Raganato, M. Apidianaki, SemEval-2024 Task 6: SHROOM, a shared-task on hallucinations and related observable overgeneration mistakes, in: *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, Association for Computational Linguistics, Mexico City, Mexico, 2024.
- [7] W. Zhao, T. Goyal, Y. Y. Chiu, L. Jiang, B. Newman, A. Ravichander, K. Chandu, R. L. Bras, C. Cardie, Y. Deng, et al., WildHallucinations: Evaluating Long-form Factuality in LLMs with Real-World Entity Queries, *arXiv preprint arXiv:2407.17468* (2024).
- [8] S. Harnad, Language Writ Large: LLMs, ChatGPT, Grounding, Meaning and Understanding, *arXiv preprint arXiv:2402.02243* (2024).
- [9] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, H. Wang, Retrieval-augmented generation for large language models: A survey, *arXiv preprint arXiv:2312.10997* (2023).
- [10] L. Hu, Z. Liu, Z. Zhao, L. Hou, L. Nie, J. Li, A Survey of Knowledge Enhanced Pre-Trained Language Models, *IEEE Transactions on Knowledge and Data Engineering* 36 (2024) 1413–1430.
- [11] F. Petroni, T. Rocktäschel, P. Lewis, A. Bakhtin, Y. Wu, A. H. Miller, S. Riedel, Language models as knowledge bases?, *arXiv preprint arXiv:1909.01066* (2019).
- [12] J.-C. Kalo, L. Fichtel, KAMEL: Knowledge Analysis with Multitoken Entities in Language Models., in: *AKBC*, 2022.
- [13] S. Min, K. Krishna, X. Lyu, M. Lewis, W.-t. Yih, P. W. Koh, M. Iyyer, L. Zettlemoyer, H. Hajishirzi, FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation, *arXiv preprint arXiv:2305.14251* (2023).
- [14] D. Plunkett, T. Sundell, Varieties of metalinguistic negotiation, *Topoi* 42 (2023) 983–999.
- [15] D. Plunkett, T. Sundell, Disagreement and the semantics of normative and evaluative terms, *Philosophers* 13 (2013).
- [16] R. E. Rudolph, Contested metalinguistic negotiation, *Synthese* 202 (2023) 90.
- [17] H. Paulheim, Knowledge graph refinement: A survey of approaches and evaluation methods, *Semantic web* 8 (2017) 489–508.
- [18] B. P. Allen, Conceptual Engineering Using Large Language Models, *arXiv preprint arXiv:2312.03749* (2023). *arXiv:2312.03749*.
- [19] V. Khatrri, C. V. Brown, Designing data governance, *Communications of the ACM* 53 (2010) 148–152.
- [20] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L. B. da Silva Santos, P. E. Bourne, et al., The FAIR Guiding Principles for scientific data management and stewardship, *Scientific data* 3 (2016) 1–9.
- [21] L. Vogt, P. Strömert, N. Matentzoglou, N. Karam, M. Konrad, M. Prinz, R. Baum, FAIR 2.0: Extending the FAIR Guiding Principles to Address Semantic Interoperability, *arXiv preprint arXiv:2405.03345* (2024).
- [22] H. Elsahar, P. Vougiouklis, A. Remaci, C. Gravier, J. Hare, F. Laforest, E. Simperl, T-REx: A



- Large Scale Alignment of Natural Language with Knowledge Base Triples, in: Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), 2018.
- [23] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, Y. Iwasawa, Large language models are zero-shot reasoners, *Advances in neural information processing systems* 35 (2022) 22199–22213.
  - [24] B. P. Allen, P. T. Groth, Evaluating Class Membership Relations in Knowledge Graphs using Large Language Models, in: *European Semantic Web Conference*, 2024. [arXiv:arXiv:2404.17000](#), to appear.
  - [25] C.-H. Chiang, H.-y. Lee, Can large language models be an alternative to human evaluations?, *arXiv preprint arXiv:2305.01937* (2023).
  - [26] A. Bavaresco, R. Bernardi, L. Bertolazzi, D. Elliott, R. Fernández, A. Gatt, E. Ghaleb, M. Giulianelli, M. Hanna, A. Koller, et al., LLMs instead of human judges? a large scale empirical study across 20 NLP evaluation tasks, *arXiv preprint arXiv:2406.18403* (2024).
  - [27] A. S. Thakur, K. Choudhary, V. S. Ramayapally, S. Vaidyanathan, D. Hupkes, Judging the Judges: Evaluating Alignment and Vulnerabilities in LLMs-as-Judges, *arXiv preprint arXiv:2406.12624* (2024).
  - [28] B. Allen, F. Polat, P. Groth, SHROOM-INDElab at SemEval-2024 Task 6: Zero- and Few-Shot LLM-Based Classification for Hallucination Detection, in: *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, Association for Computational Linguistics, Mexico City, Mexico, 2024. URL: <https://doi.org/10.48550/arXiv.2404.03732>.
  - [29] E. M. Bender, A. Koller, Climbing towards NLU: On meaning, form, and understanding in the age of data, in: *Proceedings of the 58th annual meeting of the association for computational linguistics*, 2020, pp. 5185–5198.
  - [30] M. Mandelkern, T. Linzen, Do Language Models’ Words Refer?, *arXiv preprint arXiv:2308.05576* (2024). [arXiv:2308.05576](#).
  - [31] H. Lederman, K. Mahowald, Are Language Models More Like Libraries or Like Librarians? Bibliotechnism, the Novel Reference Problem, and the Attitudes of LLMs, *arXiv preprint arXiv:2401.04854* (2024). [arXiv:2401.04854](#).
  - [32] B. A. Levinstein, D. A. Herrmann, Still no lie detector for language models: Probing empirical and conceptual roadblocks, *Philosophical Studies* (2024) 1–27.
  - [33] G. Baggio, E. Murphy, On the referential capacity of language models: An internalist rejoinder to Mandelkern & Linzen, *arXiv preprint arXiv:2406.00159* (2024).
  - [34] J. Grindrod, Large language models and linguistic intentionality, *Synthese* 204 (2024) 71.
  - [35] F. Berto, *Topics of thought: The logic of knowledge, belief, imagination*, Oxford University Press, 2022.
  - [36] P. Hawke, Theories of aboutness, *Australasian Journal of Philosophy* 96 (2018) 697–723.
  - [37] P. Hawke, L. Hornischer, F. Berto, Truth, topicality, and transparency: one-component versus two-component semantics, *Linguistics and Philosophy* (2024) 1–23.
  - [38] D. A. Herrmann, B. A. Levinstein, Standards for Belief Representations in LLMs, *arXiv preprint arXiv:2405.21030* (2024).
  - [39] J. Harding, Operationalising representation in natural language processing, *The British Journal for the Philosophy of Science* (2023). To appear.