

# Rethinking Robustness Evaluation for Question Answering: From Synthetic Stress Tests to Natural Language Variation

Anonymous ACL submission

## Abstract

In this position paper, we contend that prevailing robustness evaluation practices for Question Answering (QA) do not adequately capture system behavior under real-world conditions. Current evaluations predominantly rely on synthetic perturbations defined by idealized assumptions about linguistic validity and label preservation, whose relevance to deployment scenarios is often unclear. Consequently, robustness assessed in such settings may provide a distorted view of the reliability of QA systems. This limitation becomes increasingly salient as Large Language Models (LLMs) are deployed in interactive and agent-based applications, where language variation emerges organically and compounds across multiple interactions. We examine commonly adopted synthetic perturbation paradigms, analyze their limitations, and contrast them with emerging efforts that evaluate robustness using naturally occurring perturbations. Building on this analysis, we advocate for a community-wide shift toward robustness evaluation grounded in real-world language variation and more reliable evaluation protocols.

## 1 Introduction

Despite strong benchmark performance on Question Answering (QA) (Yang et al., 2025), Large Language Models (LLMs) can often exhibit brittle behaviour when faced with variations in textual input (Zhang et al., 2025). These failures highlight a fundamental challenge for deploying LLM-based QA systems in real-world environments, where language is inherently variable, noisy, and continuously evolving (Markov et al., 2023; Kumar and Mishra, 2025).

Robustness evaluation has therefore become a central concern. The dominant paradigm relies on *synthetic perturbations*, which apply predefined input transformations such as adversarial sentence insertion (Jia and Liang, 2017; Tran et al., 2023),

character swaps (Fang et al., 2023), or paraphrasing (Wu et al., 2023), to stress-test QA systems. More recently, a smaller but growing body of work has shifted attention toward *natural perturbations* that arise organically from real-world language variation (Wu et al., 2025b,a), rather than from artificially imposed input manipulations.

In this position paper, we argue that **current QA robustness evaluation paradigms are misaligned with how LLMs fail in practice, and that meaningful progress requires a shift from synthetic, assumption-driven perturbations toward evaluations grounded in naturally occurring language variation**. While synthetic perturbations offer control diagnostics and yield valuable insights into model behaviour (Wang et al., 2022b; Ho et al., 2023), their relevance to real-world deployment remains unclear, raising questions about their ecological validity as practical robustness probes (Kumar and Mishra, 2025). Moreover, synthetic approaches typically rely on strong assumptions regarding label preservation, human answerability, and linguistic validity; in practice, however, these assumptions are often violated by the perturbations themselves, thereby undermining the reliability of robustness evaluation (Dyrmishi et al., 2023). In contrast, natural perturbations better capture real-world conditions and tend to preserve adversarial validity more reliably (Le et al., 2022; Wu et al., 2025b); yet, they remain underexplored and rarely integrated into existing robustness evaluation protocols.

The paper first examines dominant synthetic perturbation approaches used for QA, tracing their application from early language models to modern LLMs, and analyses their limitations. It then presents natural perturbations that emerge from naturally occurring linguistic variation in practice. Based on this analysis, we outline concrete directions for rethinking robustness evaluation protocols for LLM-based QA.

## 2 Synthetic Perturbations: What Have We Learned?

A significant body of work on NLP robustness assessment relies on synthetic perturbations—deliberate modifications based on predefined input transformation strategies. These methods assume that the gold label is either preserved or altered under bounded perturbations. Accordingly, synthetic perturbation techniques can be broadly classified into two categories: label-preserving and label-changing. In the following, we briefly trace the trajectory of synthetic perturbation-based robustness evaluation work in a representative task—Question Answering. For a more detailed and comprehensive survey on robustness evaluation in QA and broader coverage across other NLP tasks, we refer readers to Ho et al. (2023) and Wang et al. (2022b); Schlegel et al. (2023), respectively.

**Label-preserving** The majority of existing work adopts the label-preserving<sup>1</sup> assumption, employing synthetic textual perturbations like the insertion of adversarial distracting sentence (Jia and Liang, 2017; Wang and Bansal, 2018; Chen et al., 2022a; Tran et al., 2023), the rephrasing of the question (Gan and Ng, 2019) or the reading paragraph (Wu et al., 2021, 2023), the addition of misinformation (Pan et al., 2023) and character-level manipulations such as character swaps (Si et al., 2021).

**Label-changing** Another line of work introduces small but meaningful input perturbations that intentionally alter the ground truth label, with the expectation that the model should adapt its prediction to reflect the change (Gardner et al., 2020; Schlegel et al., 2021; Geva et al., 2022).

Synthetic perturbations introduced in earlier work have primarily been applied to pre-LLM neural QA models. A consistent finding is that, despite achieving strong, human-comparable performance on held-out test sets (Devlin et al., 2019; Liu et al., 2019; He et al., 2021), these models exhibit varying degrees of performance degradation under synthetic perturbation settings—highlighting their reliance on statistical shortcuts to bypass genuine task requirements and exposing a lack of robustness (Ho et al., 2023).

More recently, state-of-the-art (SOTA) LLMs have demonstrated superior natural language understanding across a wide range of NLP

<sup>1</sup>Note that label-preserving does not necessarily imply semantics-preserving.

tasks—including QA (OpenAI et al., 2024; OLMo et al., 2025; Yang et al., 2025); they also outperform their fine-tuned pre-trained language model predecessors in terms of robustness (Fang et al., 2023). Nevertheless, synthetic perturbation approaches developed for earlier QA systems remain relevant, as LLMs have also been shown to be vulnerable to such perturbations. For instance, Levy et al. (2023) adapted the adversarial distracting sentence injection technique originally proposed in (Jia and Liang, 2017), prompting a strong GPT-4 model to generate a distractor sentence that answers the question similar to the original but with one critical detail changed—referred to as the “almost detail”. The instruction encouraged the model to reuse much of the original question’s phrasing while omitting the actual answer. This perturbation strategy was later found to successfully mislead less proficient GPT-TURBO, GPT3.5, and even GPT-4 itself. Fang et al. (2023) empirically investigated the effects of diverse synthetic perturbations (e.g., neighboring character swaps, synonym replacements, and combinations of multiple attack methodologies) on other LLMs such as LLAMA (Touvron et al., 2023a), and observed similar patterns of robustness failure. Besides, recent work by Bhuiya et al. (2024) revealed that, despite not requiring downstream task-specific fine-tuning like earlier QA systems, leading LLMs including GPT, LLAMA 2 (Touvron et al., 2023b) and MIXTRAL 8x7B (Jiang et al., 2024) still tend to exploit simplifying cues to circumvent the requirement to perform multi-hop reasoning. This is evidenced by their poor generalisability under a controlled challenge setting, in which distractor paragraphs were introduced to present seemingly plausible yet incorrect alternative reasoning paths, while ensuring that the correct final answer remained unchanged.

## 3 Natural Perturbations: Reflecting Real-World Failure Modes

Unlike synthetic perturbations, natural perturbations originate from authentic variations observed in real-world scenarios and are therefore considered more relevant for evaluating real-world robustness (Wu et al., 2025b). Note that the term “natural” is overloaded in NLP literature, where it can also refer either to the extent to which synthetically modified text preserves linguistic characteristics such as fluency, coherence, grammaticality, and clarity, i.e., its naturalness (Jin et al., 2020; Li et al.,

2020; Schlegel et al., 2021; Qi et al., 2021; Wang et al., 2022a; Dyrnishi et al., 2023), or to naturally occurring out-of-distribution data shift (Wang et al., 2022b). This contrasts with our focus, where “natural” pertains to perturbations arising from real-world scenarios rather than those engineered artificially. Some works also propose that a natural synthetically perturbed sample should be imperceptible to human judges (Li et al., 2020; Garg and Ramakrishnan, 2020) or convey the impression of human authorship (Dyrnishi et al., 2023). However, this proposition remains a subject of debate (Zhao et al., 2018; Wang et al., 2022b; Chen et al., 2022b). Table 1 summarises a non-exhaustive body of literature on natural perturbation methods applied across diverse NLP tasks. In the following, we categorise these works by the sources of natural perturbations.

**Wikipedia edit histories** Wikipedia’s revision histories provide a rich source of human-authored textual changes over time, offering a valuable corpus for studying real-world text variations. As one of the earlier efforts, Belinkov and Bisk (2018) explored robustness in Neural Machine Translation (NMT) by applying single-word perturbations to non-English source-side sentences. They built a lookup table of lexical errors, such as typos and misspellings, extracted from French (Max and Wisniewski, 2010) and German (Zesch, 2012) Wikipedia edit histories. Words in the source sentences were then replaced with corresponding errors from the table, where applicable. Eger and Benz (2020) extended the same approach to POS tagging, Natural Language Inference (NLI), and Toxic Comment Classification (TC), leveraging revision histories from English Wikipedia. Building on this line of research, natural perturbations were further generalised to various Question Answering (QA) tasks (Wu et al., 2025b). Instead of perturbing individual tokens, they substituted entire reading paragraphs with their edited counterparts retrieved from English Wikipedia revision histories, enabling evaluation under naturally occurring, context-level perturbations. The study assessed the sensitivity of neural language models ranging from early BERT-based architectures to SOTA LLMs, revealing that robustness issues persist across model generations.

**Human-written text** Textual content authored by human writers may contain a diverse range of errors and thus serve as a potential source of natural perturbations, as exemplified by essays written by

non-native Czech speakers (Šebesta et al., 2017; Belinkov and Bisk, 2018) and by more than 18 million sentences produced by internet users across nine real-life datasets (Le et al., 2022).

**Human-guided input variations** Some studies involve recruiting human annotators to manually craft or verify input variations. We categorise such variations as natural when annotators are not informed that their modifications will be used to fool the model. In this setting, edits are guided by the annotators’ own judgments of necessity and tend to reflect real-world scenarios, rather than being made with the explicit goal of inducing model failure (Wallace et al., 2019; Bartolo et al., 2020)—thus preserving their naturalness. Sun et al. (2024) examined the robustness of instruction-tuned models to instruction rephrasing across more than 80 NLP tasks drawn from MMLU (Hendrycks et al., 2021) and BIG-BENCH LITE (Srivastava et al., 2023). A total of 36 NLP graduate researchers were recruited to write novel instructions they believed would best elicit the desired behavior for each task. These newly crafted (unobserved) instruction phrasings, while differing superficially from those seen during instruction fine-tuning, were shown to consistently degrade model performance—highlighting limitations in the models’ generalisability. Rather than having human annotators directly propose perturbations, Chen et al. (2025) introduced 21 types of real-world scenario variations targeting natural language descriptions in LLM-based code generation tasks, derived from survey responses collected from professionals in industry and research institutions with experience using LLMs for code generation. Similarly, interview responses from employees at 16 British, German, and American NGOs, whose work directly involves online hate, were used to design 29 real-world functional tests aimed at revealing specific weaknesses in hate speech detection models (Röttger et al., 2021).

#### 4 Call for Action

We call on the NLP community to reposition robustness evaluation around naturally occurring language variation as the primary testbed for assessing LLM reliability. Concretely, this entails building and maintaining evaluation benchmarks derived from real-world linguistic evolution (e.g., edits, revisions, and user-driven reformulations), treating synthetic perturbations as auxiliary diagnostics rather than the main evaluation signal. Robust-

Task/Reference	Natural Perturbation Method	Level	Validity	Defense Strategies
NMT (Belinkov and Bisk, 2018)	Replace each word in the source-side sentence with available edits mined from French and German <i>Wikipedia edit histories</i> and <i>human-written essays</i> in Czech	Word	Unclear	Average character embedding; Training on perturbed data
POS tagging, NLI, TC (Eger and Benz, 2020)	Replace words with natural human errors from the <i>Wikipedia edit history</i>	Word	Unclear	Training on perturbed data
QA (Wu et al., 2025b)	Replace the reading passage with its counterparts based on available English <i>Wikipedia edit histories</i>	Context	Human	Adversarial training with perturbed data; In-context demonstrations
Toxic Comments/Hate Speech/Online Cyberbullying Texts Detection (Le et al., 2022)	Retrieve and substitute words using perturbations extracted from <i>a large corpus of over 18M sentences written by netizens</i> , based on phonetic similarity and edit distance	Word	Human	Sound-Invariant CNN; Adversarial training with perturbed data
Over 80 unique tasks from MMLU and BIG-BENCH LITE (Sun et al., 2024)	Recruit <i>36 NLP graduate students</i> to compose novel instructions that are appropriate for a given task but superficially different from those seen during instruction fine-tuning	Instruction	Unclear	Aligning representations of equivalent instructions
LLMs Code Generation (Chen et al., 2025)	Apply perturbations from 21 specific categories to the original natural language prompt. These categories, which may occur in real-world scenarios, were suggested by <i>experienced practitioners from the open-source community, industry, and academia</i> through the online survey	Mix	Human	–

Table 1: A non-exhaustive summary of existing literature on natural perturbations. For each work, we list the studied task, method for generating naturally perturbed test data (with *italicized underlined text* indicating the corpus source), and the perturbation level. We also indicate whether the work verifies the validity of the adversarial examples and whether it proposes defense strategies. “Unclear” denotes that no systematic experiments were conducted, though some works discuss or qualitatively assess validity.

ness studies should explicitly verify perturbation validity at scale, reporting not only performance degradation but also rates of invalid, ambiguous, or unanswerable cases. As LLMs are increasingly deployed as interactive agents, robustness evaluation should further account for perturbation accumulation across multi-step interactions, where small natural variations can compound into systemic fail-

ures. Finally, the community should adopt evaluation protocols that explicitly control for benchmark leakage and prioritize reliability-oriented metrics such as consistency, calibration, and recovery behavior over single-point accuracy. Together, these actions aim to align robustness evaluation with the conditions under which LLM-based systems operate in practice.

## 299 Limitations

300 Our analysis focuses primarily on robustness evalu-  
301 ation for QA and related agentic applications. The  
302 extent to which our conclusions generalize to other  
303 NLP tasks, such as generation-heavy or multimodal  
304 settings, requires further investigation and is be-  
305 yond the scope of this position paper. While we  
306 argue for prioritizing naturally occurring perturba-  
307 tions, constructing large-scale, high-quality natural  
308 perturbation benchmarks poses significant practical  
309 challenges, including data collection, annotation  
310 cost, and potential domain bias. We therefore call  
311 for further discussion and research on these chal-  
312 lenges.

## 313 References

314 Max Bartolo, Alastair Roberts, Johannes Welbl, Sebas-  
315 tian Riedel, and Pontus Stenetorp. 2020. [Beat the AI:  
316 Investigating adversarial human annotation for read-  
317 ing comprehension](#). *Transactions of the Association  
318 for Computational Linguistics*, 8:662–678.

319 Yonatan Belinkov and Yonatan Bisk. 2018. [Synthetic  
320 and natural noise both break neural machine transla-  
321 tion](#). In *International Conference on Learning Rep-  
322 resentations*.

323 Neeladri Bhuiya, Viktor Schlegel, and Stefan Winkler.  
324 2024. [Seemingly plausible distractors in multi-hop  
325 reasoning: Are large language models attentive read-  
326 ers?](#) In *Proceedings of the 2024 Conference on  
327 Empirical Methods in Natural Language Processing*,  
328 pages 2514–2528, Miami, Florida, USA. Association  
329 for Computational Linguistics.

330 Howard Chen, Jacqueline He, Karthik Narasimhan, and  
331 Danqi Chen. 2022a. [Can rationalization improve ro-  
332 bustness?](#) In *Proceedings of the 2022 Conference of  
333 the North American Chapter of the Association for  
334 Computational Linguistics: Human Language Tech-  
335 nologies*, pages 3792–3805, Seattle, United States.  
336 Association for Computational Linguistics.

337 Junkai Chen, Li Zhenhao, Hu Xing, and Xia Xin. 2025.  
338 [Nlperturbator: Studying the robustness of code llms  
339 to natural language variations](#). *ACM Trans. Softw.  
340 Eng. Methodol.* Just Accepted.

341 Yangyi Chen, Hongcheng Gao, Ganqu Cui, Fanchao  
342 Qi, Longtao Huang, Zhiyuan Liu, and Maosong Sun.  
343 2022b. [Why should adversarial perturbations be im-  
344 perceptible? rethink the research paradigm in adver-  
345 sarial NLP](#). In *Proceedings of the 2022 Conference  
346 on Empirical Methods in Natural Language Process-  
347 ing*, pages 11222–11237, Abu Dhabi, United Arab  
348 Emirates. Association for Computational Linguistics.

349 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
350 Kristina Toutanova. 2019. [BERT: Pre-training of](#)

[deep bidirectional transformers for language under-  
standing](#). In *Proceedings of the 2019 Conference of  
the North American Chapter of the Association for  
Computational Linguistics: Human Language Tech-  
nologies, Volume 1 (Long and Short Papers)*, pages  
4171–4186, Minneapolis, Minnesota. Association for  
Computational Linguistics.

358 Salijona Dyrnishi, Salah Ghamizi, and Maxime Cordy.  
359 2023. [How do humans perceive adversarial text?  
360 a reality check on the validity and naturalness of  
361 word-based adversarial attacks](#). In *Proceedings of the  
362 61st Annual Meeting of the Association for Compu-  
363 tational Linguistics (Volume 1: Long Papers)*, pages  
364 8822–8836, Toronto, Canada. Association for Com-  
365 putational Linguistics.

366 Steffen Eger and Yannik Benz. 2020. [From hero to  
367 zéro: A benchmark of low-level adversarial attacks](#).  
368 In *Proceedings of the 1st Conference of the Asia-  
369 Pacific Chapter of the Association for Computational  
370 Linguistics and the 10th International Joint Confer-  
371 ence on Natural Language Processing*, pages 786–  
372 803, Suzhou, China. Association for Computational  
373 Linguistics.

374 Jingliang Fang, Hua Xu, Zhijing Wu, Kai Gao, Xiaoyin  
375 Che, and Haotian Hui. 2023. [Robustness-eva-mrc:  
376 Assessing and analyzing the robustness of neural  
377 models in extractive machine reading comprehension](#).  
378 *Intelligent Systems with Applications*, 20:200287.

379 Wee Chung Gan and Hwee Tou Ng. 2019. [Improv-  
380 ing the robustness of question answering systems to  
381 question paraphrasing](#). In *Proceedings of the 57th  
382 Annual Meeting of the Association for Computational  
383 Linguistics*, pages 6065–6075, Florence, Italy. Asso-  
384 ciation for Computational Linguistics.

385 Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan  
386 Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi,  
387 Dheeru Dua, Yanai Elazar, Ananth Gottumukkala,  
388 Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco,  
389 Daniel Khashabi, Kevin Lin, Jiangming Liu, Nel-  
390 son F. Liu, Phoebe Mulcaire, Qiang Ning, and 7 oth-  
391 ers. 2020. [Evaluating models’ local decision bound-  
392 aries via contrast sets](#). In *Findings of the Association  
393 for Computational Linguistics: EMNLP 2020*, pages  
394 1307–1323, Online. Association for Computational  
395 Linguistics.

396 Siddhant Garg and Goutham Ramakrishnan. 2020.  
397 [BAE: BERT-based adversarial examples for text clas-  
398 sification](#). In *Proceedings of the 2020 Conference on  
399 Empirical Methods in Natural Language Processing  
400 (EMNLP)*, pages 6174–6181, Online. Association for  
401 Computational Linguistics.

402 Mor Geva, Tomer Wolfson, and Jonathan Berant. 2022.  
403 [Break, perturb, build: Automatic perturbation of rea-  
404 soning paths through question decomposition](#). *Trans-  
405 actions of the Association for Computational Linguis-  
406 tics*, 10:111–126.



523	<a href="#">for evaluating comprehension in machine reading.</a>	<a href="#">comprehension models.</a>	580
524	<i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 35(15):13762–13770.	In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 1543–1557, Dubrovnik, Croatia. Association for Computational Linguistics.	581
525			582
526	Viktor Schlegel, Goran Nenadic, and Riza Batista-Navarro. 2023. <a href="#">A survey of methods for revealing and overcoming weaknesses of data-driven natural language understanding.</a>		583
527	<i>Natural Language Engineering</i> , 29(1):1–31.		584
528		Eric Wallace, Pedro Rodriguez, Shi Feng, Ikuya Yamada, and Jordan Boyd-Graber. 2019. <a href="#">Trick me if you can: Human-in-the-loop generation of adversarial examples for question answering.</a>	585
529		<i>Transactions of the Association for Computational Linguistics</i> , 7:387–401.	586
530			587
531	Karel Šebesta, Zuzanna Bedřichová, Kateřina Šormová, Barbora Štindlová, Milan Hrdlička, Tereza Hrdličková, Jiří Hana, Vladimír Petkevič, Tomáš Jelínek, Svatava Škodová, Petr Janeš, Kateřina Lundáková, Hana Skoumalová, Šimon Sládek, Piotr Pierscieniak, Dagmar Toufarová, Milan Straka, Alexandr Rosen, Jakub Náplava, and Marie Poláčková. 2017. <a href="#">CzeSL grammatical error correction dataset (CzeSL-GEC).</a>		588
532	LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL).		589
533		Jiayi Wang, Rongzhou Bao, Zhuosheng Zhang, and Hai Zhao. 2022a. <a href="#">Distinguishing non-natural from natural adversarial samples for more robust pre-trained language model.</a>	590
534		In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 905–915, Dublin, Ireland. Association for Computational Linguistics.	591
535			592
536		Xuezhi Wang, Haohan Wang, and Diyi Yang. 2022b. <a href="#">Measure and improve robustness in NLP models: A survey.</a>	593
537		In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4569–4586, Seattle, United States. Association for Computational Linguistics.	594
538			595
539			596
540			597
541			598
542	Chenglei Si, Ziqing Yang, Yiming Cui, Wentao Ma, Ting Liu, and Shijin Wang. 2021. <a href="#">Benchmarking robustness of machine reading comprehension models.</a>		599
543	In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 634–644, Online. Association for Computational Linguistics.		600
544			601
545			602
546			603
547			604
548	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazary, and 431 others. 2023. <a href="#">Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.</a>		605
549	<i>Transactions on Machine Learning Research</i> . Featured Certification.		606
550		Yicheng Wang and Mohit Bansal. 2018. <a href="#">Robust machine comprehension models via adversarial training.</a>	607
551		In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 575–581, New Orleans, Louisiana. Association for Computational Linguistics.	608
552			609
553			610
554			611
555			612
556		Winston Wu, Dustin Arendt, and Svitlana Volkova. 2021. <a href="#">Evaluating neural model robustness for machine comprehension.</a>	613
557		In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2470–2481, Online. Association for Computational Linguistics.	614
558	Jiuding Sun, Chantal Shaib, and Byron C Wallace. 2024. <a href="#">Evaluating the zero-shot robustness of instruction-tuned language models.</a>		615
559	In <i>The Twelfth International Conference on Learning Representations</i> .		616
560			617
561			618
562			619
563	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. <a href="#">Llama: Open and efficient foundation language models.</a>		620
564	<i>Preprint</i> , arXiv:2302.13971.	Yulong Wu, Viktor Schlegel, and Riza Batista-Navarro. 2023. <a href="#">Are machine reading comprehension systems robust to context paraphrasing?</a>	621
565		In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 184–196, Nusa Dua, Bali. Association for Computational Linguistics.	622
566			623
567			624
568			625
569			626
570	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrubhi Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023b. <a href="#">Llama 2: Open foundation and fine-tuned chat models.</a>		627
571	<i>Preprint</i> , arXiv:2307.09288.		628
572		Yulong Wu, Viktor Schlegel, and Riza Batista-Navarro. 2025a. <a href="#">Natural context drift undermines the natural language understanding of large language models.</a>	629
573		In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 1248–1259, Suzhou, China. Association for Computational Linguistics.	630
574			631
575			632
576			633
577	Son Quoc Tran, Phong Nguyen-Thuan Do, Uyen Le, and Matt Kretzmar. 2023. <a href="#">The impacts of unanswerable questions on the robustness of machine reading</a>		634
578		Yulong Wu, Viktor Schlegel, and Riza Batista-Navarro. 2025b. <a href="#">Pay attention to real world perturbations! nat-</a>	635
579			636

637 [ural robustness evaluation in machine reading com-](#)  
638 [prehension](#). *Preprint*, arXiv:2502.16523.

639 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,  
640 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,  
641 Chengen Huang, Chenxu Lv, Chujie Zheng, Day-  
642 iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao  
643 Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41  
644 others. 2025. [Qwen3 technical report](#). *Preprint*,  
645 arXiv:2505.09388.

646 Torsten Zesch. 2012. [Measuring contextual fitness using](#)  
647 [error contexts extracted from the Wikipedia re-](#)  
648 [vision history](#). In *Proceedings of the 13th Confer-*  
649 *ence of the European Chapter of the Association for*  
650 *Computational Linguistics*, pages 529–538, Avignon,  
651 France. Association for Computational Linguistics.

652 Kun Zhang, Le Wu, Kui Yu, Guangyi Lv, and Dacao  
653 Zhang. 2025. [Evaluating and improving robustness](#)  
654 [in large language models: A survey and future direc-](#)  
655 [tions](#). *Preprint*, arXiv:2506.11111.

656 Zhengli Zhao, Dheeru Dua, and Sameer Singh. 2018.  
657 [Generating natural adversarial examples](#). In *Internat-*  
658 *ional Conference on Learning Representations*.