

Pay Attention to Real World Perturbations! Natural Robustness Evaluation in Machine Reading Comprehension

Anonymous ACL submission

Abstract

As neural language models achieve human-comparable performance on Machine Reading Comprehension (MRC) and see widespread adoption, ensuring their robustness in real-world scenarios has become increasingly important. Current robustness evaluation research, though, primarily develops synthetic perturbation methods, leaving unclear how well they reflect real life. Considering this, we present a framework to automatically examine MRC models on occurring textual perturbations, by replacing paragraph in MRC benchmarks with their counterparts based on available Wikipedia edit history. Such perturbation type is *natural* as its design does not stem from an artificial generative process, inherently distinct from the previously investigated synthetic approaches. In a large-scale study encompassing various model architectures we observe that natural perturbations result in performance degradation in pre-trained encoder language models, with errors extending to Flan-T5 and Large Language Models (LLMs). We also show that exposing encoder-only models to naturally perturbed examples during training contributes to handling natural perturbations. This adversarial training approach, however, is not able to promote performance improvement on the majority of synthetic perturbations, indicating that many types of synthetic noise do not actually exist in our collected real-world textual perturbations. We hope this study will inspire future robustness investigation efforts to focus more on natural perturbations, thus deepening our understanding of how models respond to realistic linguistic challenges and providing insights into practical robustness enhancement strategies.

1 Introduction

Transformer-based pre-trained language models demonstrate remarkable efficacy in addressing questions based on a given passage of text, a task commonly referred to as Machine Reading Comprehension (MRC) (Devlin et al., 2019; Brown

et al., 2020; He et al., 2021; Wei et al., 2022; Tournon et al., 2023; OpenAI et al., 2024). Despite these advancements, high-performing MRC systems are also known to succeed by relying on *shortcuts* in benchmark datasets rather than truly demonstrating understanding of the passage, thereby lacking robustness to various types of test-time perturbations (Ho et al., 2023; Schlegel et al., 2023; Levy et al., 2023).

Evaluating models' resilience to textual perturbations during inference aids in identifying adversarial instances that highlight their shortcut behavior and provides insights into mitigating these shortcuts. While numerous synthetic perturbation approaches have been explored and reveal the vulnerabilities of MRC models to various linguistic challenges (Ribeiro et al., 2018; Jiang and Bansal, 2019; Welbl et al., 2020; Schlegel et al., 2021; Cao et al., 2022; Tran et al., 2023), a serious concern is that these carefully designed perturbations might not necessarily appear in real-world settings. Consequently, this poses a risk of neglecting the weaknesses of reading comprehension systems to real challenges when deployed in practical scenarios, thus potentially hindering the improvement of their reliability in practical applications.

To counteract this issue, in this paper, we develop a framework to inject textual changes that arise in real-world conditions into MRC datasets and audit how well contemporary language models perform under such perturbations. We deem them as natural because the perturbation process does not involve any artificial manipulation, in line with the definitions by Belinkov and Bisk (2018); Hendrycks et al. (2021); Pedraza et al. (2022); Agarwal et al. (2022). Results of robustness evaluation are therefore more representative of real-world applications. Inspired by Belinkov and Bisk (2018), our approach utilises English Wikipedia revision histories as the source of natural perturbations, given that the differences between revi-

sions authentically capture the textual modifications made by human editors in the real world. By comparing the variances between each adjacent revision, we identify perturbed versions for each Wikipedia reading passage in the original MRC benchmarks (if it exists). Apart from natural perturbations, we also investigate the robustness of MRC models to diverse synthetic perturbation techniques at different levels to discern their differences. All perturbation methods only alter the reading context, while the questions and ground truth answers remain unchanged.

With the established framework, we conduct extensive experiments on two datasets, evaluating twenty-nine models, including nine recently proposed LLMs. Experimental results indicate that natural perturbations encompass rich linguistic variations and can lead to failures in the encoder-only models, while humans are almost undeterred by their presence. Crucially, these errors also transfer to larger and more powerful models, such as FLan-T5 and state-of-the-art LLMs. On synthetic perturbations, we also observe a lack of robustness across all model architectures, although the realism of thusly generated adversarial examples remains a concern. Using naturally perturbed MRC instances for retraining usually boosts the robustness of encoder-only models against natural perturbations. However, this often comes with a decline in original performance. Retraining also sometimes helps improve robustness to synthetic perturbations, but for the majority of cases, it does not enhance performance on the perturbed test set and even decrease it, likely due to the inherent differences between natural and synthetic perturbations.

2 Related Work

Robustness Evaluation in MRC A typical approach to evaluate the robustness of MRC models is via test-time perturbation. This line of research develops different perturbation methods as attacks, such as adversarial distracting sentence addition (Jia and Liang, 2017; Tran et al., 2023), word substitution (Wu et al., 2021), character swap (Si et al., 2021), entity renaming (Yan et al., 2022) and paraphrasing (Gan and Ng, 2019; Lai et al., 2021; Wu et al., 2023a). Our work also fits within the category of test-time perturbation, but differs from previous works in that we introduce perturbations that naturally occur in real-world scenarios, therefore contributing to a more practical robustness exami-

nation. We also experiment with various synthetic perturbations for comparison purposes.

Natural Perturbation for Robustness Assessment Compared with deliberately crafting the perturbed instances, the study of natural perturbation is under-explored. In the computer vision domain, researchers find that real-world clean images without intentional modifications can confuse deep learning models as well, terming them as natural adversarial examples (Hendrycks et al., 2021; Pedraza et al., 2022). Similarly, in the field of Natural language processing (NLP), Belinkov and Bisk (2018) concludes that naturally occurring errors dramatically break machine translation systems. Motivated by these, we attempt to harvest natural perturbations from available Wikipedia revision histories and utilise them to modify the original MRC instances. To the best of our knowledge, we are the first to investigate MRC model robustness under real natural perturbations. Furthermore, it should be noted that the concept of natural perturbed examples in this paper differs from what is defined in previous NLP literature, where the latter measures the extent to which synthetically modified text preserves certain 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). 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).

3 Natural and Synthetic Perturbation

In this section, we detail our methodology to create label-preserving stress MRC test sets by introducing real-world occurring noises and artificial perturbations to the context paragraph, respectively. The size of each examined MRC dataset and the correspondingly constructed perturbed test sets are displayed in Appendix A.

3.1 Natural Perturbation

We design a pipeline to automatically construct MRC test sets with noises that occur in real-world settings by leveraging Wikipedia revision histories. Our approach comprises two modules: *candidate*

184 *passage pairs curation and perturbed test set con-*
185 *struction.*

186 **Candidate passage pairs curation.** For each En-
187 glish Wikipedia article within the development set
188 of MRC datasets, we systematically extract its en-
189 tire revision histories and preprocess them, includ-
190 ing the removal of markups and the segmentation
191 of content. Subsequently, we obtain the content
192 differences between each current revision and the
193 previous adjacent one, identifying three distinct
194 editing patterns: addition, deletion, and modifi-
195 cation¹. In the case of an edit falling within the
196 modification pattern, we retain the paragraph from
197 the prior version as the *original* and the correspond-
198 ing one from the current version as the *perturbed*,
199 provided both paragraphs exceed 500 characters².
200 This results in a total of 91,093 pairs of candidate
201 reading passages distributed across 46 articles.

202 **Perturbed test set construction.** To generate the
203 naturally perturbed test set, we begin by acquir-
204 ing all reading passages from the development set
205 of each MRC dataset and identifying their entries
206 in the collection of previously extracted candidate
207 original passages, along with the corresponding per-
208 turbed counterparts. Subsequently, for the matched
209 original passages with a single occurrence, we keep
210 them and the corresponding perturbed passages;
211 whereas for those with multiple occurrences, we
212 randomly select one instance for each and extract
213 its perturbed version. After obtaining the perturbed
214 reading passages, we retain only those with at least
215 one question where all annotated ground truth an-
216 swers (or all plausible answers for the unanswer-
217 able question) can still be located within the per-
218 turbed context, resulting in the *Perturbed* test set.
219 For the sake of comparison, we also construct an
220 *Original* version of the test set keeping only the
221 original passages and questions corresponding to
222 those that were included in the *Perturbed* version.

223 3.2 Synthetic Perturbation

224 To explore the difference between natural and syn-
225 thetic perturbations, our study incorporates a com-
226 prehensive range of synthetic perturbation tech-
227 niques, spanning various linguistic levels: char-
228 acter, word, sentence and document-level, as de-

¹In Appendix B, the average percentage of editing patterns observed for each Wikipedia article in the investigated MRC development datasets is presented, revealing that “modification” constitutes the predominant editing pattern.

²This threshold setting adheres to the methodology employed in the collection of SQuAD 1.1 (Rajpurkar et al., 2016).

229 tailed in Table 1. While certain character-level
230 and word-level perturbation methods have been in-
231 vestigated across multiple NLP tasks, such as the
232 CharSwapMid for machine translation (Belinkov
233 and Bisk, 2018) and the WDelete for quality es-
234 timation (Kanojia et al., 2021; Wu et al., 2023b),
235 none of these has been applied to the contextual
236 paragraph to study the robustness for the task of ex-
237 tractive MRC. We employ methods including WS-
238 plit, WSynSub and WInsert (WE) to each sentence
239 in the original reading passage, and then recombine
240 the modified sentences to generate the perturbed
241 version. Conversely, other perturbation approaches
242 are directly executed on the entire paragraph, as im-
243 plementing them at the sentence-level might result
244 in perturbed text that is even difficult for humans
245 to read and comprehend (Si et al., 2021). The im-
246 plementation of all character-level and word-level
247 methods is carried out using the NLPAug library
248 (Ma, 2019). Moreover, we set the perturbation rate
249 to 30%, in line with the default settings within the
250 NLPAug library.

251 To conduct a comprehensive evaluation, we
252 also apply a variation of the sentence-level per-
253 turbation method known as AddSent (Jia and
254 Liang, 2017; Chen et al., 2022a; Tran et al., 2023;
255 Levy et al., 2023). Our approach prompts the
256 GPT-3.5-turbo-0125 model to generate a distrac-
257 tor sentence that shares significant lexical overlap
258 with the question but is not an appropriate answer
259 to it and is also irrelevant to the context. Unlike pre-
260 vious approaches, we do not explicitly require the
261 distractor sentence to be a statement that answers
262 the so-called “almost detail” question (Levy et al.,
263 2023). Instead, the model can generate any type
264 of sentence as long as it satisfies the three criteria
265 stated in the prompt. Subsequently, from the set
266 of generated candidate distractor sentence-question
267 pairs under each reading passage, we select the one
268 with the highest lexical overlap. We finally insert
269 the distractor sentence from the identified pair at
270 the beginning of the original context, as previous lit-
271 erature suggests that prepending results in a larger
272 impact on the performance (Ko et al., 2020; Chen
273 et al., 2022a). For document-level perturbations,
274 we introduce two methods that both leverage the
275 capabilities of the GPT-3.5-turbo-0125 model as
276 well. The former DocPara attempts to directly para-
277 phrase the entire context paragraph, while the latter
278 Style Transfer, drawing inspiration from (Qi et al.,
279 2021), seeks to transfer the style of the reading
280 passage by rephrasing it using a distinct persona

discerned based on its topic. The manually constructed prompt for methods AddSent, DocPara and Style Transfer are shown in Appendix C.

For each perturbation method, excluding AddSent, which preserves all ground truth answers in nature, we conduct the same answers-preserving checking as described in Section 3.1 after obtaining the perturbed reading passages to construct the corresponding *Original* and *Perturbed* test set pair for evaluation purposes.

Method	Explanation
<i>character-level</i>	
CharOCR	Replace characters with predefined Optical Character Recognition (OCR) errors.
CharInsert	Inject new characters randomly.
CharSubstitute	Substitute original characters randomly.
CharSwapMid	Swap adjacent characters within words randomly, excluding the first and last character.
CharSwapRand	Swap characters randomly without constraint.
<i>word-level</i>	
WInsert (CWE)	Insert new words to random position according to contextual word embeddings calculation from RoBERTa-base.
WSubstitute (CWE)	Substitute words according to contextual word embeddings calculation from RoBERTa-base (Liu et al., 2019).
WSplit	Split words to two tokens randomly.
WSwap	Swap adjacent words randomly.
WDelete	Delete words randomly.
WCrop	Remove a set of continuous word randomly.
Word Synonym Substitution (WSynSub)	Substitute words with synonyms from large size English PPDB (Pavlick et al., 2015).
WInsert (WE)	Insert new words to random position according to GloVe (Pennington et al., 2014) word embeddings calculation ⁵ .
<i>sentence-level</i>	
AddSent	Add a context-irrelevant distractor sentence with high lexical overlap to the question at the beginning of the context.
<i>document-level</i>	
Document Paraphrasing (DocPara)	Paraphrasing the whole context paragraph directly.
Style Transfer	Rephrase the passage using a distinct persona discerned based on its topic.

Table 1: Various synthetic perturbation approaches.

4 Experiments Setup

4.1 Datasets

We select two widely studied benchmark MRC datasets (License: CC-BY-SA-4.0) for which human performance has been surpassed by state-of-the-art models, due to the fact that their reading passages are sourced from Wikipedia, thereby enabling the utilisation of Wikipedia editing histories to generate the naturally perturbed test set.

SQuAD 1.1 (Rajpurkar et al., 2016): An English reading comprehension dataset with over 100,000 questions created by crowdworkers on a set of Wikipedia article paragraphs. Each question is accompanied by multiple ground truth answers, each of which represents a continuous span from the corresponding reading passage.

SQuAD 2.0 (Rajpurkar et al., 2018): The combination of SQuAD 1.1 with over 50,000 unanswerable questions crafted adversarially by crowdworkers, thus considered to be more challenging.

4.2 Models

Our evaluation study involves multiple contemporary MRC models across three different types: encoder-only, encoder-decoder, and decoder-only. Under the encoder-decoder and decoder-only model evaluation settings, we reframe the extractive MRC as the text generation task based on the given context and question. Access to and experimentation with all models are possible via the use of the HuggingFace’s *Transformers* library (Wolf et al., 2020), two 80GB Nvidia A100 GPUs and the OpenAI ChatGPT API.

Encoder-only: We select BERT (Devlin et al., 2019) and its various variants for evaluation, including DistilBERT (Sanh et al., 2019), SpanBERT (Joshi et al., 2020), RoBERTa (Liu et al., 2019), ALBERT (Lan et al., 2020) and DeBERTa (He et al., 2021). Some of these model types also come with different variations, such as size (e.g., *base* and *large* for RoBERTa), versions (e.g., *v1* and *v2* for ALBERT) and whether the input text is cased or not (e.g., *cased* and *uncased* for BERT), all of which are included in the evaluation. We fine-tune these encoder-only pre-trained language models on the training set of our examined MRC datasets and evaluate them on the constructed original and perturbed test sets. Model details and the hyperparameters used in model fine-tuning are shown in Appendix D.

Encoder-Decoder: Instruction finetuning has been demonstrated to be effective in enhancing zero-shot performance of pretrained language models, resulting in the development of Finetuned Language Net (FLAN) (Wei et al., 2022). In this work, we use the instruction-finetuned version of T5 model class, specifically the FLan-T5 (Chung et al., 2022), available in sizes ranging from *small* (80M), *base* (250M), *large* (780M) to *xl* (3B). During evaluation, we utilise the instruction templates from MRC task collection in open-sourced FLAN repository and report the model performance as the average of those obtained across the employed templates. Refer to Appendix E for various instruction templates used for the evaluation on the test sets with the format as the two examined MRC datasets.

Decoder-only: There is an exponential increase of pre-trained generative LLMs and their fine-tuned chat versions, inspired by the remarkable success of ChatGPT (Bang et al., 2023). Therefore, our experiments incorporate a broad range of recently proposed language model families, including GPT 3.5 Turbo, Llama 2 (Touvron et al., 2023), Llama 3, Mistral (Jiang et al., 2023), Falcon (Almazrouei et al., 2023) and Gemma (Mesnard et al., 2024). The zero-shot prompts designed for soliciting responses from them are presented in Appendix F.

4.3 Model Evaluation Metrics

We choose the Exact Match (EM) and (Macro-averaged) F1 score to assess the performance of both encoder-only and encoder-decoder models, as they strictly output the shortest continuous span from the context as the answer (or predict the question as unanswerable) during inference. However, for almost all of the decoder-only models, their outputs are not consistently adhere to the instruction due to their conversational style, rendering EM and F1 metrics unsuitable for evaluation. Consequently, we employ a more lenient metric, namely Inclusion Match (IM), which measures whether the response of the model contains any of the ground truth answers. Furthermore, if the model’s output includes phrases such as “I cannot answer this/the question” and “unanswerable”⁴, we deem that the model believes the question is not answerable. Model robustness is quantified by measuring the relative variation in performance (as reflected in the F1 or IM) under perturbations.

5 MRC under Natural Perturbation

5.1 Are Encoder-only MRC Models Resilient to Natural Perturbation?

Table 2 presents the relative F1 change for all encoder-only MRC models on the naturally perturbed test set generated based on the SQuAD 1.1 and SQuAD 2.0 development set, respectively. It can be clearly seen from Table 2 that overall, the performance of all the examined models decreases, indicating that *encoder-only MRC models suffer from natural perturbation*. However, we notice that the performance drop of all models is negligible (the biggest drop is only 3.06%), which suggests that those models also exhibit considerable

⁴We identify a collection of such phrases by manually examining the decoder-only models’ outputs (Check Appendix G for the full set).

robustness to natural perturbations.

Dataset	Attacker	
	Victim	Nature
SQuAD 1.1	distilbert-base	-0.6
	bert-base-cased	-0.21
	bert-base-uncased	-0.87
	bert-large-cased	-0.63
	bert-large-uncased	-0.35
	spanbert-base-cased	-0.26
	spanbert-large-cased	-0.51
	roberta-base	-0.61
	roberta-large	-0.29
	albert-base-v1	-1.0
	albert-base-v2	-0.34
	albert-large-v1	-0.42
	albert-large-v2	-0.8
	albert-xxlarge-v1	-0.75
albert-xxlarge-v2	-0.46	
deberta-large	-0.52	
SQuAD 2.0	distilbert-base	-0.71 (-2.76/1.71)
	bert-base-cased	-0.63 (-1.84/0.6)
	bert-base-uncased	-0.49 (-1.88/0.94)
	bert-large-cased	-0.53 (-1.61/0.55)
	bert-large-uncased	-1.38 (-2.51/-0.24)
	spanbert-base-cased	-1.24 (-2.66/0.15)
	spanbert-large-cased	-1.2 (-1.9/-0.56)
	roberta-base	-0.6 (-2.09/0.81)
	roberta-large	-1.52 (-2.6/-0.54)
	albert-base-v1	-1.07 (-2.02/-0.22)
	albert-base-v2	-1.08 (-2.03/-0.22)
	albert-large-v1	-0.41 (-1.42/0.52)
	albert-large-v2	-0.69 (-1.66/0.22)
	albert-xxlarge-v1	-1.23 (-3.06/0.49)
albert-xxlarge-v2	-1.28 (-3.02/0.36)	
deberta-large	-1.05 (-2.2/0.0)	

Table 2: Relative F1 change (%) for encoder-only MRC systems subjecting to natural perturbations. In SQuAD 2.0, the values shown in the parentheses represent the relative change for answerable and unanswerable questions, respectively.

5.2 Error Analysis

Although encoder-only MRC models exhibit a relatively small performance gap, it remains worthwhile to investigate the sources of natural perturbation and reveal the perturbation phenomena contributing to models’ error. To this end, within the original and the perturbed test set pair generated using the “Nature” method based on SQuAD 2.0 development set, we first identify 384 instances where at least one encoder-only model succeeds on the original but fails⁵ on the perturbed (i.e., being adversarial), and then randomly select the same num-

⁵For answerable questions, a model’s prediction is considered correct if both the EM and F1 scores are 1, and incorrect if both metrics are 0 or it determines the question is unanswerable. For unanswerable questions, a model’s prediction is correct if it predicts the question is unanswerable, and wrong if it provides an answer span.

ber of instances on which all encoder-only models succeed on both the original and perturbed versions (Naik et al., 2018). We refer to these two types of instances as C2W (correct to wrong) and C2C (correct to correct) instances, respectively. Among the identified C2W and C2C instances, we further remove duplicates, resulting in 210 and 244 unique original and perturbed paragraph pairs, respectively. Furthermore, as natural perturbation can occasionally help the model to get the answer correct, we also filter 85 unique W2C (wrong to correct) instances on which at least two encoder-only models fail on the original but succeed on the perturbed. Finally, utilising an 8-category taxonomy of the semantic edit intentions in Wikipedia revisions derived from Yang et al. (2017), the first author of the paper manually annotated the chosen 210 samples of C2W and C2C, as well as the 85 W2C samples. To validate our findings, we further present 20% of the annotated C2W and C2C examples to a second annotator. See Appendix H for the instruction provided to the annotators, along with detailed explanations of each edit intention. We calculate the (micro-averaged) F1 score to evaluate the inter-annotator agreement, which is 0.82. This suggests that the annotators’ annotations align closely. Table 3 reports the annotation results.

Edit Intention	C2W	C2C	W2C
Copy Editing	43.3	47.1	40.0
Clarification	5.7	3.3	1.2
Elaboration	23.8	18.1	22.4
Fact Update	4.3	3.8	3.5
Refactoring	1.9	1.9	1.2
Simplification	14.3	8.6	21.2
Vandalism	21.0	17.6	23.5
Other	9.5	14.8	10.6

Table 3: The percentage (%) of samples annotated with each edit intention in the C2W, C2C and W2C categories. The percentages do not add up to 100% because a single revision may fall into multiple intentions.

From Table 3, we observe that there is no significant difference in the distribution of annotated edit intentions between C2W and C2C examples, suggesting that *though these types of natural perturbations confuse the encoder-only MRC models, the effect is not as pronounced*. A roughly similar distribution is also observed in the W2C examples, which indicates that these natural perturbation types can also facilitate correct answers

by the models, i.e., being beneficial. Copy editing constitutes the most frequent edit intention (more than 40%), followed by elaboration and vandalism, with refactoring represents the category with the lowest percentage. Moreover, we find that there might be no correlation between the quality of the perturbed passage and its potential for being adversarial in the MRC robustness evaluation. Certain text edits aimed at improving the passage quality, such as copy editing and elaboration, do render the perturbation adversarial, whereas edits intended to damage the article may not consistently result in adversarial instances; in fact, vandalism can even assist models in providing correct answers. Instead, we infer that whether an edit to the passage can render the MRC instance adversarial or not depends on the location of the edits in relation to the question. Among the 384 C2W and C2C examples, we measure the proportion of answerable questions with the answer sentence(s) in the original passage remaining unmodified in the naturally perturbed version, which is 34.5% and 71.5%, respectively. This confirms our hypothesis that if the edits affect the answer sentence(s), there is a higher likelihood of the perturbed example becoming adversarial; otherwise, it might not. Appendix I presents one perturbed example for each of the C2W, C2C, and W2C categories, respectively, along with the annotated natural perturbation type(s).

5.3 Validity of Nature Adversarial Examples

To accurately assess a model’s robustness under perturbation, it is vital to examine the validity of adversarial example, i.e. whether humans can still find the correct answer under the perturbation (Dyrmishi et al., 2023). We first present two human annotators with the same collection of adversarial instances, which includes only perturbed contexts and their corresponding questions, and then ask them to answer the question based on the perturbed context. The annotators are required to select the shortest continuous span in the perturbed context that answers the question(s) and are allowed to leave the answer blank if they are confident that the question is not answerable. Full instructions given to the annotators can be seen in Appendix H. Subsequently, for both annotators, we measure the correctness (1 or 0) of their provided answers by comparing each of them with the corresponding ground truth answers⁶. The inter-annotator agree-

⁶Here, as long as one of the ground truth answers is included in the human-provided answer span, we consider the

ment is then measured by computing the Cohen’s κ coefficient (Cohen, 1960). We then involve a third human annotator to annotate the adversarial examples on which the first two annotators disagree and then take the majority label as ground truth.

We employ this approach to verify the validity of the 210 C2W examples in Section 5.2 and find that 86% of these adversarial examples are valid (0.77 Cohen’s κ), indicating that *a substantial proportion of natural adversarial examples for encoder-only MRC model(s) are valid.*

5.4 Can Errors from Encoder-only Models Affect Other Architectures?

We are also curious about how well the errors identified in encoder-only models carry over to other model architectures. This leads us to first propose an exhaustive search algorithm that leverages the predictions of all encoder-only models to create the challenging natural perturbed test set. In detailed terms, for each matched reading passage from the prior version and its counterpart from the current version, we determine which should be designated as the *original* and which as the *perturbed* based on which scenario can yield the questions on which the maximum sum of the number of encoder-only models demonstrates the lack of robustness phenomenon⁷. Questions on which none of the encoder-only models fail under the perturbation are then removed. We finally process the identified original and perturbed passage pairs to ensure that the original passages are within the original SQuAD 1.1 development set. For those original passages with multiple occurrences, we select the one with the maximum number of questions reserved. With the development set of SQuAD 1.1 and SQuAD 2.0 as the source, this results in two challenge perturbed test sets: NAT_V1_CHALLENGE and NAT_V2_CHALLENGE. In NAT_V1_CHALLENGE, there are 184 contexts and 234 questions. NAT_V2_CHALLENGE contains 214 contexts and 442 questions (226 unanswerable).

Table 4 shows the evaluation results of both encoder-decoder and decoder-only models on the newly generated challenge test sets. From the table, we observe that *the errors caused by natu-*

⁷A model lacks robustness to the perturbation if it achieves 1 EM on the original but attains less than 0.4 F1 on the perturbed (for answerable questions).

ral perturbation in encoder-only MRC models transfer to both Flan-T5 and LLMs. On the NAT_V1_CHALLENGE, Flan-T5-small demonstrates the greatest susceptibility to natural perturbation, experiencing a 14.27% decrease in F1, while among LLMs, Gemma-7B-IT emerges as the least robust, with a 16.66% F1 drop. Transitioning to the NAT_V2_CHALLENGE, the base version of Flan-T5 exhibits the largest performance decline (13.83%) and Falcon-7B-Instruct stands out as the LLM with the lowest robustness. In Appendix J, we showcase two adversarial examples targeting LLMs sourced from our generated challenge sets.

Model	Performance <i>original vs. perturbed</i>			
	NAT_V1_CHALLENGE		NAT_V2_CHALLENGE	
	original	perturbed	original	perturbed
flan-t5-small	58.76/64.76	48.58/55.52 _{-14.27}	42.57/44.57	39.71/41.81 _{-6.19}
flan-t5-base	79.49/85.01	66.17/73.42 _{-13.63}	70.66/72.85	61.16/62.78 _{-13.83}
flan-t5-large	88.1/92.53	76.57/82.31 _{-11.05}	79.11/81.01	70.14/72.13 _{-10.96}
flan-t5-xl	86.25/91.57	75.0/81.45 _{-11.05}	83.71/85.84	73.19/74.86 _{-12.79}
GPT-3.5-turbo-0125	91.03	83.33 _{-8.46}	51.58	47.06 _{-8.76}
Gemma-2B-IT	51.28	43.16 _{-15.83}	55.66	50.23 _{-9.76}
Gemma-7B-IT	82.05	68.38 _{-16.66}	59.95	57.01 _{-4.9}
Llama 2-chat-7B	82.91	73.93 _{-10.83}	41.63	38.69 _{-7.06}
Llama 2-chat-13B	80.77	73.93 _{-8.47}	46.83	41.18 _{-12.06}
Llama-3-8B-Instruct	88.89	77.35 _{-12.98}	51.81	46.61 _{-10.04}
Mistral-7B-Instruct-v0.2	85.9	76.92 _{-10.45}	55.43	52.04 _{-6.12}
Falcon-7B-Instruct	53.42	50 _{-6.4}	32.81	23.53 _{-28.28}
Falcon-40B-Instruct	69.66	62.82 _{-9.82}	38.69	36.88 _{-4.68}

Table 4: The performance (%) of encoder-decoder and decoder-only MRC models on the newly generated original and naturally perturbed challenge test sets. Values in smaller font are changes (%) relative to the original performance of the model.

6 MRC Under Synthetic Perturbation

In Appendix K, we present the evaluation results of different levels of synthetic perturbations against all MRC model architectures on the correspondingly generated test dataset. It can be seen from this table that generally, *MRC systems exhibit a lack of robustness to synthetic perturbations, with varying degrees of performance decline.* Methods AddSent, WSplit and WInsert (WE) lead to noticeable drops in model performance, whereas other techniques demonstrate relatively limited impact.

For each method, from its created SQuAD 2.0-format test set pair, we also randomly select 50 instances where the GPT-3.5-turbo-0125 shows evidence of being not robust, resulting in a total of 800 adversarial examples. We then measure their validity using the methodology described in Section 5.3, shuffling their order to mitigate potential bias, and present the results (0.81 Cohen’s κ) in Table 5. From Table 5, we can see that in general, character and word-level perturbation methods result in more valid adversarial

examples than sentence and document-level approaches, even though certain methods are lacking validity, such as WSubstitute (CWE) (48%) and CharSwapRand (52%). WSplit achieves the highest attack validity with 74%, while AddSent attains the lowest with 28%, despite causing the largest performance decrease for the GPT-3.5-turbo-0125. This suggests that the AddSent method, while impactful, might frequently generate perturbed MRC instances on which even humans find challenging. In Appendix L, we demonstrate some valid synthetic adversarial examples.

Attack	Answered Correctly
CharOCR	64
CharInsert	70
CharSubstitute	56
CharSwapMid	60
CharSwapRand	52
WInsert (CWE)	64
WSubstitute (CWE)	48
WSplit	74
WSwap	60
WDelete	60
WCrop	68
WSynSub	62
WInsert (WE)	58
AddSent	28
DocPara	48
Style Transfer	46

Table 5: The percentage (%) of adversarial MRC instances correctly labelled by humans for each synthetic perturbation method.

7 Adversarial Training

To enhance model robustness, we conduct adversarial training by identifying six encoder-only model architectures that exhibit the highest robustness to natural perturbations in their respective categories (except albert-xxlarge-v2 on NAT_V2_CHALLENGE), and presenting them with both original training data and the generated naturally perturbed training examples. We extract the entire Wikipedia revision histories for the 392 articles in the original SQuAD training set, and then obtain 5,262 (with 22,033 questions) and 5,311 (with 32,993 questions) perturbed contexts to augment the original SQuAD 1.1 and SQuAD 2.0 training set, respectively, using the methodology described in Section 3.1. Table 6 compares the performance of these models on NAT_V1_CHALLENGE and NAT_V2_CHALLENGE, before and after retraining. Further, we also evaluate the behavior of the retrained models on the constructed syn-

thetically perturbed test sets and quantify the discrepancy from the performance achieved prior to retraining. The results are shown in Appendix M.

Model	Performance			
	<i>original vs. perturbed</i>			
	NAT_V1_CHALLENGE		NAT_V2_CHALLENGE	
distilbert-base	64.53/70.45	41.03/47.6 _{-32.43}	56.56/59.08	41.18/43.3 _{-26.71}
	57.26/63.44	43.59/51.87 _{-18.24}	53.17/55.4	43.89/45.51 _{-17.85}
bert-large-cased	79.06/83.66	63.68/70.23 _{-16.05}	66.29/68.35	53.17/55.04 _{-19.47}
	74.79/80.14	59.83/67.5 _{-15.77}	67.87/69.31	58.37/59.53 _{-14.11}
spanbert-large-cased	84.19/88.2	67.95/74.77 _{-15.23}	78.73/80.68	62.44/64.99 _{-19.45}
	82.48/86.6	69.66/76.05 _{-12.18}	78.28/80.0	65.61/67.12 _{-16.1}
roberta-large	86.75/90.21	73.93/79.47 _{-11.91}	82.13/84.27	66.29/68.52 _{-18.69}
	83.33/87.15	70.94/76.53 _{-12.19}	81.22/82.67	70.59/71.84 _{-13.1}
albert-xxlarge-v2	84.62/89.64	73.93/78.77 _{-12.13}	84.62/86.07	68.11/69.61 _{-19.12}
	86.32/90.93	75.64/81.07 _{-10.84}	82.58/84.08	70.59/72.78 _{-13.44}
deberta-large	88.46/92.5	73.57/78.48 _{-15.16}	85.07/86.65	71.49/73.0 _{-15.75}
	88.03/91.84	76.92/81.53 _{-11.23}	83.03/85.1	72.62/74.48 _{-12.48}

Table 6: Comparison of the performance of several encoder-only MRC systems on NAT_V1_CHALLENGE and NAT_V2_CHALLENGE, before and after retraining. The results shown in the shaded areas represent the performance of the model retrained on the augmented training set with naturally perturbed instances.

Overall, we observe that retraining enhances both the performance of the models on the naturally perturbed test set and their robustness to natural perturbations, albeit causing a slight decrease in the original performance. However, the phenomenon of improved perturbed performance does not generally apply to most synthetic perturbations, which indicates that natural and synthetic perturbations might indeed be different.

8 Conclusion

In this paper, we mainly study the robustness of MRC models to natural perturbations, which occur under real-world conditions without intentional human intervention. Using the proposed evaluation framework, we show that certain naturally perturbed examples can indeed be adversarial, i.e., lead to model failure, even when the modifications aim to improve the overall passage quality. Natural perturbations also appear to differ significantly from synthetic ones, exhibiting a wide range of rich linguistic phenomena and may be more effective in generating valid adversarial instances. Adversarial training via augmentation with naturally perturbed samples is generally beneficial for enhancing the model’s robustness to natural perturbations; yet, it is not particularly successful in handling most synthetic noises. Future work includes the exploration of alternative natural perturbation approaches and the design of more effective defensive strategies against both natural and synthetic attacks.

647 Limitations

648 We acknowledge the presence of several limita-
649 tions in our work: (i) There is a need to expand
650 our study to other MRC datasets to make the find-
651 ings more generalisable; (ii) Our natural perturba-
652 tion framework only works with Wikipedia-based
653 benchmarks. Therefore, it is necessary to develop
654 other methods that can introduce real-world textual
655 perturbations; (iii) It is essential to design better
656 prompts to enhance the alignment of certain LLM
657 outputs with given instructions (particularly for
658 unanswerable questions), thereby ensuring more
659 accurate evaluation results. There is also a neces-
660 sity to examine the robustness of LLMs using dif-
661 ferent prompting strategies such as few-shot in-
662 context learning; (iv) Since the impact of training
663 data augmentation is relatively limited, we need to
664 explore better techniques to improve the robustness
665 of encoder-only models to natural perturbations
666 and further investigate the robustness connection
667 between natural and synthetic perturbations. En-
668 hancing the robustness of LLMs is also a potential
669 future direction.

670 References

671 Akshay Agarwal, Nalini Ratha, Mayank Vatsa, and
672 Richa Singh. 2022. Exploring robustness connec-
673 tion between artificial and natural adversarial exam-
674 ples. In *Proceedings of the IEEE/CVF Conference*
675 *on Computer Vision and Pattern Recognition (CVPR)*
676 *Workshops*, pages 179–186.

677 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Al-
678 shamsi, Alessandro Cappelli, Ruxandra Cojocaru,
679 Mérouane Debbah, Étienne Goffinet, Daniel Hesslow,
680 Julien Launay, Quentin Malartic, Daniele Mazzotta,
681 Badreddine Noune, Baptiste Pannier, and Guilherme
682 Penedo. 2023. *The falcon series of open language*
683 *models*. Preprint, arXiv:2311.16867.

684 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-
685 liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei
686 Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu,
687 and Pascale Fung. 2023. *A multitask, multilingual,*
688 *multimodal evaluation of ChatGPT on reasoning, hal-*
689 *lucination, and interactivity*. In *Proceedings of the*
690 *13th International Joint Conference on Natural Lan-*
691 *guage Processing and the 3rd Conference of the Asia-*
692 *Pacific Chapter of the Association for Computational*
693 *Linguistics (Volume 1: Long Papers)*, pages 675–718,
694 Nusa Dua, Bali. Association for Computational Lin-
695 guistics.

696 Yonatan Belinkov and Yonatan Bisk. 2018. *Synthetic*
697 *and natural noise both break neural machine transla-*
698 *tion*. In *International Conference on Learning Rep-*
699 *resentations*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie
700 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
701 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
702 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
703 Gretchen Krueger, Tom Henighan, Rewon Child,
704 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
705 Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-
706 teusz Litwin, Scott Gray, Benjamin Chess, Jack
707 Clark, Christopher Berner, Sam McCandlish, Alec
708 Radford, Ilya Sutskever, and Dario Amodei. 2020.
709 *Language models are few-shot learners*. In *Ad-*
710 *vances in Neural Information Processing Systems*,
711 volume 33, pages 1877–1901. Curran Associates,
712 Inc. 713

Yu Cao, Dianqi Li, Meng Fang, Tianyi Zhou, Jun Gao,
714 Yibing Zhan, and Dacheng Tao. 2022. *TASA: De-*
715 *ceiving question answering models by twin answer*
716 *sentences attack*. In *Proceedings of the 2022 Con-*
717 *ference on Empirical Methods in Natural Language*
718 *Processing*, pages 11975–11992, Abu Dhabi, United
719 Arab Emirates. Association for Computational Lin-
720 guistics. 721

Howard Chen, Jacqueline He, Karthik Narasimhan, and
722 Danqi Chen. 2022a. *Can rationalization improve ro-*
723 *bustness?* In *Proceedings of the 2022 Conference of*
724 *the North American Chapter of the Association for*
725 *Computational Linguistics: Human Language Tech-*
726 *nologies*, pages 3792–3805, Seattle, United States.
727 Association for Computational Linguistics. 728

Yangyi Chen, Hongcheng Gao, Ganqu Cui, Fanchao
729 Qi, Longtao Huang, Zhiyuan Liu, and Maosong Sun.
730 2022b. *Why should adversarial perturbations be im-*
731 *perceptible? rethink the research paradigm in adver-*
732 *sarial NLP*. In *Proceedings of the 2022 Conference*
733 *on Empirical Methods in Natural Language Process-*
734 *ing*, pages 11222–11237, Abu Dhabi, United Arab
735 Emirates. Association for Computational Linguistics. 736

Hyung Won Chung, Le Hou, Shayne Longpre, Barret
737 Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi
738 Wang, Mostafa Dehghani, Siddhartha Brahma, Al-
739 bert Webson, Shixiang Shane Gu, Zhuyun Dai,
740 Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-
741 ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,
742 Dasha Valter, Sharan Narang, Gaurav Mishra, Adams
743 Yu, Vincent Zhao, Yanping Huang, Andrew Dai,
744 Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-
745 cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,
746 and Jason Wei. 2022. *Scaling instruction-finetuned*
747 *language models*. Preprint, arXiv:2210.11416. 748

Jacob Cohen. 1960. *A coefficient of agreement for*
749 *nominal scales*. *Educational and Psychological Mea-*
750 *surement*, 20(1):37–46. 751

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
752 Kristina Toutanova. 2019. *BERT: Pre-training of*
753 *deep bidirectional transformers for language under-*
754 *standing*. In *Proceedings of the 2019 Conference of*
755 *the North American Chapter of the Association for*
756 *Computational Linguistics: Human Language Tech-*
757 *nologies, Volume 1 (Long and Short Papers)*, pages
758

759	4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	816
760		817
761	Salijona Dyrnishi, Salah Ghamizi, and Maxime Cordy.	818
762	2023. How do humans perceive adversarial text? a reality check on the validity and naturalness of word-based adversarial attacks. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8822–8836, Toronto, Canada. Association for Computational Linguistics.	819
763		820
764		
765		821
766		822
767		823
768		824
769		825
770	Wee Chung Gan and Hwee Tou Ng. 2019. Improving the robustness of question answering systems to question paraphrasing. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6065–6075, Florence, Italy. Association for Computational Linguistics.	826
771		827
772		828
773		829
774		830
775	Siddhant Garg and Goutham Ramakrishnan. 2020. BAE: BERT-based adversarial examples for text classification. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6174–6181, Online. Association for Computational Linguistics.	831
776		832
777		
778		833
779		834
780		835
781	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. {DEBERTA}: {DECODING}-{enhanced} {bert} {with} {disentangled} {attention}. In <i>International Conference on Learning Representations</i> .	836
782		837
783		838
784		839
785		
786	Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. 2021. Natural adversarial examples. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 15262–15271.	840
787		841
788		842
789		843
790		844
791	Xanh Ho, Johannes Mario Meissner, Saku Sugawara, and Akiko Aizawa. 2023. A survey on measuring and mitigating reasoning shortcuts in machine reading comprehension. <i>Preprint</i> , arXiv:2209.01824.	845
792		
793		846
794		847
795	Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2021–2031, Copenhagen, Denmark. Association for Computational Linguistics.	848
796		849
797		850
798		
799		851
800		852
801	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. Mistral 7b. <i>Preprint</i> , arXiv:2310.06825.	853
802		854
803		855
804		856
805		
806		857
807		858
808		859
809	Yichen Jiang and Mohit Bansal. 2019. Avoiding reasoning shortcuts: Adversarial evaluation, training, and model development for multi-hop QA. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2726–2736, Florence, Italy. Association for Computational Linguistics.	860
810		861
811		862
812		863
813		
814		864
815		865
		866
		867
		868
		869
		870
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870

871	Gemma Team Thomas Mesnard, Cassidy Hardin,	Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	933
872	Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	934
873	L. Sifre, Morgane Riviere, Mihir Kale, J Christo-	Simón Posada Fishman, Juston Forte, Isabella Ful-	935
874	pher Love, Pouya Dehghani Tafti, L'eonard Hussenot,	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	936
875	Aakanksha Chowdhery, Adam Roberts, Aditya	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	937
876	Barua, Alex Botey, Alex Castro-Ros, Ambrose	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	938
877	Slone, Am'elie H'eliou, Andrea Tacchetti, Anna Bu-	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	939
878	lanova, Antonia Paterson, Beth Tsai, Bobak Shahri-	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	940
879	ari, Charline Le Lan, Christopher A. Choquette-Choo,	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	941
880	Cl'ement Crepy, Daniel Cer, Daphne Ippolito, David	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	942
881	Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	943
882	Yan, George Tucker, George-Christian Muraru, Grig-	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	944
883	ory Rozhdestvenskiy, Henryk Michalewski, Ian Ten-	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	945
884	ney, Ivan Grishchenko, Jacob Austin, James Keel-	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	946
885	ing, Jane Labanowski, Jean-Baptiste Lespiau, Jeff	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	947
886	Stanway, Jenny Brennan, Jeremy Chen, Johan Fer-	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	948
887	ret, Justin Chiu, Justin Mao-Jones, Katherine Lee,	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	949
888	Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa	Christina Kim, Yongjik Kim, Jan Hendrik Kircher-	950
889	Lee, Lucas Dixon, Machel Reid, Maciej Mikula,	ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo,	951
890	Mateo Wirth, Michael Sharman, Nikolai Chinaev,	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	952
891	Nithum Thain, Olivier Bachem, Oscar Chang, Oscar	stantinidis, Kyle Kopic, Gretchen Krueger, Vishal	953
892	Wahltimez, Paige Bailey, Paul Michel, Petko Yotov,	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	954
893	Pier Giuseppe Sessa, Rahma Chaabouni, Ramona	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	955
894	Comanescu, Reena Jana, Rohan Anil, Ross McIl-	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	956
895	roy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Se-	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	957
896	bastian Borgeaud, Sertan Girgin, Sholto Douglas,	Anna Makanju, Kim Malfacini, Sam Manning, Todor	958
897	Shree Pandya, Siamak Shakeri, Soham De, Ted Kli-	Markov, Yaniv Markovski, Bianca Martin, Katie	959
898	menko, Tom Hennigan, Vladimir Feinberg, Woj-	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	960
899	ciech Stokowiec, Yu hui Chen, Zafarali Ahmed,	McKinney, Christine McLeavey, Paul McMillan,	961
900	Zhitao Gong, Tris Brian Warkentin, Ludovic Peran,	Jake McNeil, David Medina, Aalok Mehta, Jacob	962
901	Minh Giang, Cl'ement Farabet, Oriol Vinyals, Jeffrey	Menick, Luke Metz, Andrey Mishchenko, Pamela	963
902	Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	964
903	Ghahramani, Douglas Eck, Joelle Barral, Fernando	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	965
904	Pereira, Eli Collins, Armand Joulin, Noah Fiedel,	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	966
905	Evan Senter, Alek Andreev, and Kathleen Kenealy.	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	967
906	2024. Gemma: Open models based on gemini re-	Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex	968
907	search and technology . <i>ArXiv</i> , abs/2403.08295.	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	969
908	Aakanksha Naik, Abhilasha Ravichander, Norman	tista Parascandolo, Joel Parish, Emy Parparita, Alex	970
909	Sadeh, Carolyn Rose, and Graham Neubig. 2018.	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	971
910	Stress test evaluation for natural language inference .	man, Filipe de Avila Belbute Peres, Michael Petrov,	972
911	In <i>Proceedings of the 27th International Conference</i>	Henrique Ponde de Oliveira Pinto, Michael, Poko-	973
912	<i>on Computational Linguistics</i> , pages 2340–2353,	rnell, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	974
913	Santa Fe, New Mexico, USA. Association for Com-	erly, Alethea Power, Boris Power, Elizabeth Proehl,	975
914	putational Linguistics.	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	976
915	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	Cameron Raymond, Francis Real, Kendra Rimbach,	977
916	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	978
917	man, Diogo Almeida, Janko Altmenschmidt, Sam Alt-	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	979
918	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	Girish Sastry, Heather Schmidt, David Schnurr, John	980
919	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	Schulman, Daniel Selsam, Kyla Sheppard, Toki	981
920	ing Bao, Mohammad Bavarian, Jeff Belgium, Ir-	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	982
921	wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	983
922	Christopher Berner, Lenny Bogdonoff, Oleg Boiko,	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	984
923	Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-	Sokolowsky, Yang Song, Natalie Staudacher, Felipe	985
924	man, Tim Brooks, Miles Brundage, Kevin Button,	Petroski Such, Natalie Summers, Ilya Sutskever,	986
925	Trevor Cai, Rosie Campbell, Andrew Cann, Brittany	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	987
926	Carey, Chelsea Carlson, Rory Carmichael, Brooke	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	988
927	Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	989
928	Chen, Ruby Chen, Jason Chen, Mark Chen, Ben	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	990
929	Chess, Chester Cho, Casey Chu, Hyung Won Chung,	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	991
930	Dave Cummings, Jeremiah Currier, Yunxing Dai,	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	992
931	Cory Decareaux, Thomas Degry, Noah Deutsch,	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	993
932	Damien Deville, Arka Dhar, David Dohan, Steve	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	994
		Clemens Winter, Samuel Wolrich, Hannah Wong,	995
		Lauren Workman, Sherwin Wu, Jeff Wu, Michael	996

997	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.		
998			
999			
1000			
1001			
1002			
1003	Ellie Pavlick, Pushpendre Rastogi, Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. 2015. PPDB 2.0: Better paraphrase ranking, fine-grained entailment relations, word embeddings, and style classification . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 425–430, Beijing, China. Association for Computational Linguistics.		
1004			
1005			
1006			
1007			
1008			
1009			
1010			
1011			
1012			
1013	Anibal Pedraza, Oscar Deniz, and Gloria Bueno. 2022. Really natural adversarial examples. <i>International Journal of Machine Learning and Cybernetics</i> , 13(4):1065–1077.		
1014			
1015			
1016			
1017	Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.		
1018			
1019			
1020			
1021			
1022			
1023	Fanchao Qi, Yangyi Chen, Xurui Zhang, Mukai Li, Zhiyuan Liu, and Maosong Sun. 2021. Mind the style of text! adversarial and backdoor attacks based on text style transfer . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4569–4580, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		
1024			
1025			
1026			
1027			
1028			
1029			
1030			
1031	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for SQuAD . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 784–789, Melbourne, Australia. Association for Computational Linguistics.		
1032			
1033			
1034			
1035			
1036			
1037			
1038	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.		
1039			
1040			
1041			
1042			
1043			
1044	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Semantically equivalent adversarial rules for debugging NLP models . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 856–865, Melbourne, Australia. Association for Computational Linguistics.		
1045			
1046			
1047			
1048			
1049			
1050			
1051	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter . In <i>5th</i>		
1052			
1053			
		<i>Workshop on Energy Efficient Machine Learning and Cognitive Computing @ NeurIPS 2019</i> .	1054
			1055
		Viktor Schlegel, Goran Nenadic, and Riza Batista-Navarro. 2021. Semantics altering modifications for evaluating comprehension in machine reading . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 35(15):13762–13770.	1056
			1057
			1058
			1059
			1060
		Viktor Schlegel, Goran Nenadic, and Riza Batista-Navarro. 2023. A survey of methods for revealing and overcoming weaknesses of data-driven natural language understanding . <i>Natural Language Engineering</i> , 29(1):1–31.	1061
			1062
			1063
			1064
			1065
		Chenglei Si, Ziqing Yang, Yiming Cui, Wentao Ma, Ting Liu, and Shijin Wang. 2021. Benchmarking robustness of machine reading comprehension models . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 634–644, Online. Association for Computational Linguistics.	1066
			1067
			1068
			1069
			1070
			1071
		Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>Preprint</i> , arXiv:2307.09288.	1072
			1073
			1074
			1075
			1076
			1077
			1078
			1079
			1080
			1081
			1082
			1083
			1084
			1085
			1086
			1087
			1088
			1089
			1090
			1091
			1092
			1093
			1094
		Son Quoc Tran, Phong Nguyen-Thuan Do, Uyen Le, and Matt Kretchmar. 2023. The impacts of unanswerable questions on the robustness of machine reading comprehension models . 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.	1095
			1096
			1097
			1098
			1099
			1100
			1101
			1102
		Jiayi Wang, Rongzhou Bao, Zhuosheng Zhang, and Hai Zhao. 2022a. Distinguishing non-natural from natural adversarial samples for more robust pre-trained language model . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 905–915, Dublin, Ireland. Association for Computational Linguistics.	1103
			1104
			1105
			1106
			1107
			1108
			1109
		Xuezhi Wang, Haohan Wang, and Diyi Yang. 2022b. Measure and improve robustness in NLP models: A	1110
			1111

survey. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4569–4586, Seattle, United States. Association for Computational Linguistics.

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. *Finetuned language models are zero-shot learners*. In *International Conference on Learning Representations*.

Johannes Welbl, Pasquale Minervini, Max Bartolo, Pontus Stenetorp, and Sebastian Riedel. 2020. *Under-sensitivity in neural reading comprehension*. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1152–1165, Online. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. *Transformers: State-of-the-art natural language processing*. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Winston Wu, Dustin Arendt, and Svitlana Volkova. 2021. *Evaluating neural model robustness for machine comprehension*. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2470–2481, Online. Association for Computational Linguistics.

Yulong Wu, Viktor Schlegel, and Riza Batista-Navarro. 2023a. *Are machine reading comprehension systems robust to context paraphrasing?* In *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)*, pages 184–196, Nusa Dua, Bali. Association for Computational Linguistics.

Yulong Wu, Viktor Schlegel, Daniel Beck, and Riza Batista-Navarro. 2023b. *MMT’s submission for the WMT 2023 quality estimation shared task*. In *Proceedings of the Eighth Conference on Machine Translation*, pages 856–862, Singapore. Association for Computational Linguistics.

Jun Yan, Yang Xiao, Sagnik Mukherjee, Bill Yuchen Lin, Robin Jia, and Xiang Ren. 2022. *On the robustness of reading comprehension models to entity renaming*. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 508–520, Seattle, United States. Association for Computational Linguistics.

Diyi Yang, Aaron Halfaker, Robert Kraut, and Eduard Hovy. 2017. *Identifying semantic edit intentions from revisions in Wikipedia*. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2000–2010, Copenhagen, Denmark. Association for Computational Linguistics.

Zhengli Zhao, Dheeru Dua, and Sameer Singh. 2018. *Generating natural adversarial examples*. In *International Conference on Learning Representations*.

A Datasets Statistics

We demonstrate the number of titles, contexts, and questions contained in the studied MRC datasets and the constructed perturbed test sets in Table 7.

Method	SQuAD 1.1			SQuAD 2.0		
	titles	contexts	questions	titles	contexts	questions _(HAns/NAns)
training	442	18896	87599	442	19035	130319 _(86821/43498)
development	48	2067	10570	35	1204	11873 _(5928/5945)
Nature	44	674	2776	32	368	3174 _(1508/1666)
CharOCR	48	2007	6724	35	1194	7983 _(3808/4175)
CharInsert	48	1952	5928	35	1189	7137 _(3355/3782)
CharSubstitute	48	1956	5925	35	1181	7114 _(3371/3743)
CharSwapMid	48	1956	6053	35	1188	7257 _(3423/3834)
CharSwapRand	48	1952	5937	35	1190	7206 _(3416/3790)
WInsert (CWE)	48	2030	7622	35	1203	9037 _(4297/4740)
WSubstitute (CWE)	48	2005	6943	35	1200	8234 _(3915/4319)
WSplit	48	1352	2477	35	993	3121 _(1374/1747)
WSwap	48	1987	6257	35	1189	7442 _(3516/3926)
WDelete	48	1982	6218	35	1190	7462 _(3534/3928)
WCrop	48	2038	7459	35	1200	8809 _(4240/4569)
WSynSub	48	1791	4970	35	1136	5427 _(2546/2881)
WInsert (WE)	48	1840	4792	35	1166	5816 _(2644/3172)
AddSent	48	2067	2075	35	1204	1205 _(596/609)
DocPara	48	1968	6759	35	1188	7677 _(3615/4062)
Style Transfer	48	1968	6740	35	1196	7753 _(3647/4106)

Table 7: Dataset statistics of the SQuAD 1.1 and SQuAD 2.0, along with the respective perturbed test sets generated based on the development set of each. HAns: answerable questions; NAns: unanswerable questions.

B Visualisation of the Average Percentage of Editing Patterns

Figure 1 shows the average percentage of editing patterns identified across the articles contained within the development set of the investigated MRC datasets.

C Prompts for Sentence-Level and Document-Level Perturbations

In this section, we provide the prompts developed for implementing synthetic perturbations at both the sentence-level and document-level.

AddSent: *Generate a sentence unrelated to the context that shares significant lexical overlap with the given question but is not an answer to it.*
 \n\nQuestion: {question}\n\nContext: {context}

DocPara: *Given the context, paraphrase it as much as possible while still preserving the original*

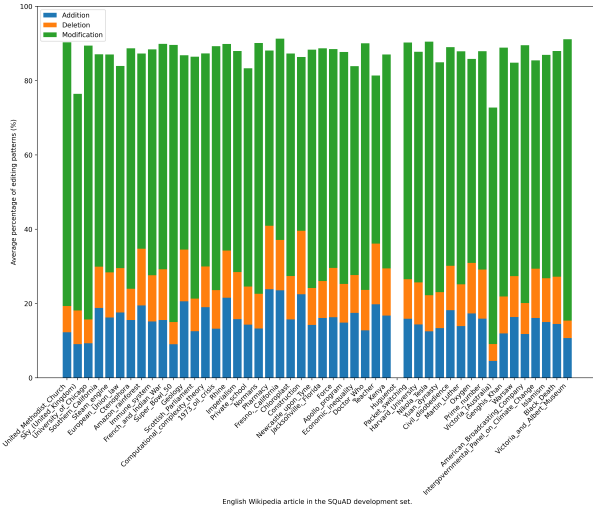


Figure 1: Average percentage of editing patterns for each Wikipedia article in the development set of our examined MRC datasets.

meaning. Make sure to keep the phrases {ground truth answers⁸} in the answer.
 Context: {context}
 The format of the output should be as follows:
 Rephrased Context:

Style Transfer: Given the context, identify a persona that can be used to rewrite it in a manner that results in a rephrased context that is both natural and realistic within a real-world scenario, while still preserving the original meaning. Make sure to keep the phrases {ground truth answers} in the answer.
 Context: {context}
 The format of the output should be as follows:
 Persona:
 Rephrased Context:

D Encoder-only Model Parameters and Hyperparameters for Fine-tuning

Table 8 shows the hyperparameters used to fine-tune the pre-trained encoder-only MRC models in this work and their number of parameters contained.

E Instruction Templates for Flan-T5 Evaluation

In Table 9, we present the instruction templates employed in constructing the inputs to the Flan-T5 model for the SQuAD 1.1 format and SQuAD 2.0 format datasets, respectively.

⁸To be precise, here we are referring to all ground truth answers annotated for the questions that belong to the given context.

Model	Parameters (M)	d	b	lr	ep
DistilBERT	(66)	384	8	3e-5	3
BERT	(110/340)	384	8	3e-5	2
SpanBERT	(110/340)	512	4	2e-5	4
RoBERTa	(125/355)	384	8	3e-5	2
ALBERT	(11/17/223)	384	4	3e-5	2
DeBERTa	(350)	384	4	3e-6	3

Table 8: Number of parameters in each type of pre-trained encoder-only MRC model and the hyperparameters used to fine-tune them. For BERT, SpanBERT, RoBERTa and ALBERT, we show the number of model parameters in the order of *base*, *large* and *xxlarge* (if applicable) version. d is the size of the token sequence fed into the model, b is the training batch size, lr is the learning rate, and ep is the number of training epochs. We used $stride = 128$ for documents longer than d tokens.

F MRC Prompts

We use the following two zero-shot prompts to instruct the decoder-only models to generate responses in the task of MRC, designed for the test sets with SQuAD 1.1 and SQuAD 2.0 format, respectively.

SQuAD 1.1: Use the provided article delimited by triple quotes to answer question. Provide only the shortest continuous span from the context without any additional explanation.
 Context: {context}
 Question: {question}

SQuAD 2.0: Use the provided article delimited by triple quotes to answer question. Provide only the shortest continuous span from the context without any additional explanation. If the question is unanswerable, return "unanswerable".
 Context: {context}
 Question: {question}

G Indicators of Unanswerable

We manually identify a set of phrases contained in the output of LLMs that indicate the unanswerability of the question, including "I cannot answer this/the question", "unanswerable", "There is no indication in the provided article", "The context provided does not provide enough information", "There is no reference in the given article", "The answer to the question is not provided in the given article", "it is not possible", "question cannot be answered" and "context/question/article/text/article provided/passage does not".

SQuAD 1.1	
1	“Read this and answer the question\n\n{context}\n\n{question}”
2	“{context}\n{question}”
3	“Answer a question about this article:\n{context}\n{question}”
4	“Here is a question about this article: {context}\nWhat is the answer to this question: {question}”
5	“Article: {context}\n\nQuestion: {question}”
6	“Article: {context}\n\nNow answer this question: {question}”
SQuAD 2.0	
1	“Read this and answer the question. If the question is unanswerable, say \"unanswerable\".\n\n{context}\n\n{question}”
2	“{context}\n{question} (If the question is unanswerable, say \"unanswerable\")”
3	“{context}\nTry to answer this question if possible (otherwise reply \"unanswerable\"): {question}”
4	“{context}\nIf it is possible to answer this question, answer it for me (else, reply \"unanswerable\"): {question}”
5	“{context}\n\nAnswer this question, if possible (if impossible, reply \"unanswerable\"): {question}”
6	“Read this: {context}\n\nNow answer this question, if there is an answer (If it cannot be answered, return \"unanswerable\"): {question}”

Table 9: Various instruction templates for the F1an-T5 model evaluation on the two benchmark MRC datasets.

H Human Annotation Instructions

In Figure 2, we show the instructions given to human annotators for error analysis (Section 5.2) and adversarial validity checking (Section 5.3), respectively. All our human annotators are students from universities in the United Kingdom and China. Before commencing each task, we ask them to annotate some examples and report the average time spent on each. As compensation, annotators receive 40 pence for each annotated example.

I Demonstration of Perturbed MRC Examples for Encoder-only Models

Figure 3 illustrates a naturally perturbed MRC instance each for categories C2W, C2C, and W2C, with the annotated perturbation type(s).

Error Analysis
<p>You will be presented with pairs of reading contexts and their modified versions. The task is to compare each context and its modified version, observe the changes made and classify them into one or more of the semantic edit intention categories detailed below:</p> <ul style="list-style-type: none"> • <i>Copy Editing</i>: Rephrase; improve grammar, spelling, tone, or punctuation • <i>Clarification</i>: Specify or explain an existing fact or meaning by example or discussion without adding new information • <i>Elaboration</i>: Extend/add new content; insert a fact or new meaningful assertion • <i>Fact Update</i>: Update numbers, dates, scores, episodes, status, etc. based on newly available information • <i>Refactoring</i>: Restructure the article; move and rewrite content, without changing the meaning of it • <i>Simplification</i>: Reduce the complexity or breadth of discussion; may remove information • <i>Vandalism</i>: Deliberately attempt to damage the article • <i>Other</i>: None of the above <p>We will use your annotation to calculate the percentage of each edit category.</p>
Adversarial Validity Checking
<p>Please read each provided context carefully and answer a corresponding question. Select the shortest continuous span from the context as your answer. If you believe a question cannot be answered, leave the answer blank. Your answer will be compared with the ground truth answers, and the result will only be used to decide the human answerability of the question.</p>

Figure 2: Instructions for the two distinct human annotation tasks. In the error analysis task, the eight semantic edit intentions are adopted from (Yang et al., 2017).

<p>Category: C2W</p> <p>Original Paragraph: <i>Jacksonville, like most large cities in the United States, suffered from negative effects of rapid urban sprawl after World War II. The construction of highways led residents to move to newer housing in the suburbs. After World War II, the government of the city of Jacksonville began to increase spending to fund new public building projects in the boom that occurred after the war. [...]</i></p> <p>Perturbed Paragraph: <i>Jacksonville, like most large cities in the United States, suffered from negative effects of rapid urban sprawl after World War V. The construction of highways led residents to move to newer housing in the suburbs. After World War II, the government of the city of Jacksonville began to increase spending to fund new public building projects in the boom that occurred after the war. [...]</i></p> <p>Question: What did Jacksonville suffer from following World War I?</p> <p>Prediction of distilbert-base and spanbert-large-cased: unanswerable→rapid urban sprawl</p> <p>Annotated Natural Perturbation Type: Vandalism</p>
<p>Category: C2C</p> <p>Original Paragraph: <i>Construction projects can suffer from preventable financial problems. Underbids happen when builders ask for too little money to complete the project. Cash flow problems exist when the present amount of funding cannot cover the current costs for labour and materials, and because they are a matter of having sufficient funds at a specific time, can arise even when the overall total is enough. Fraud is a problem in many fields, but is notoriously prevalent in the construction field. Financial planning for the project is intended to ensure that a solid plan with adequate safeguards and contingency plans are in place before the project is started and is required to ensure that the plan is properly executed over the life of the project.</i></p> <p>Perturbed Paragraph: <i>Financial planning ensures adequate safeguards and contingency plans are in place before the project is started, and ensures that the plan is properly executed over the life of the project. Construction projects can suffer from preventable financial problems. Underbids happen when builders ask for too little money to complete the project. Cash flow problems exist when the present amount of funding cannot cover the current costs for labour and materials; such problems may arise even when the overall budget is adequate, presenting a temporary issue. Fraud is also an occasional construction issue.</i></p> <p>Question: What can construction projects suffer from?</p> <p>Prediction of all encoder-only models: preventable financial problems→preventable financial problems</p> <p>Annotated Natural Perturbation Type: Copy Editing; Refactoring; Simplification</p>
<p>Category: W2C</p> <p>Original Paragraph: <i>[...] The antigens expressed by tumors have several sources; some are derived from oncogenic viruses like human papillomavirus, which causes cervical cancer, while others are the organism's own proteins that occur at low levels in normal cells but reach high levels in tumor cells. [...] A third possible source of tumor antigens are proteins normally important for regulating cell growth and survival, that commonly mutate into cancer inducing molecules called oncogenes.</i></p> <p>Perturbed Paragraph: <i>[...] The antigens expressed by tumors have several sources; some are derived from oncogenic viruses like human papillomavirus, which causes cancer of the cervix, vulva, vagina, penis, anus, mouth, and throat, while others are the organism's own proteins that occur at low levels in normal cells but reach high levels in tumor cells. [...] A third possible source of tumor antigens are proteins normally important for regulating cell growth and survival, that commonly mutate into cancer inducing molecules called oncogenes.</i></p> <p>Question: What is a fourth possible source for tumor antigens?</p> <p>Prediction of bert-base-uncased: proteins normally important for regulating cell growth and survival→unanswerable</p> <p>Annotated Natural Perturbation Type: Elaboration</p>

Figure 3: Natural perturbed MRC example in C2W, C2C and W2C categories.

1273	J Natural Adversarial Samples for LLMs	L Synthetic Adversarial Instances	1282
1274	We demonstrate two naturally perturbed reading	In Figure 5 and Figure 6, we present several syn-	1283
1275	comprehension examples that pose challenges for	thetic adversarial samples that can be solved by	1284
1276	LLMs in Figure 4.	humans.	1285
1277	K Evaluation Results Under Synthetic	M Robustness Connection Between	1286
1278	Perturbation	Synthetic and Natural Perturbations	1287
1279	The complete robustness evaluation results of the	Figure 7 describes the impact of natural adversarial	1288
1280	MRC models under synthetic perturbations are	retraining on handling synthetic perturbations and	1289
1281	shown in Table 10.	will be included in the main body of the paper once	1290
		we are given one additional page of content.	1291

NAT_V1_CHALLENGE

Original Paragraph: *In business, notable alumni include Microsoft CEO Satya Nadella, Oracle Corporation founder and the third richest man in America Larry Ellison, Goldman Sachs and MF Global CEO as well as former Governor of New Jersey Jon Corzine, McKinsey & Company founder and author of the first management accounting textbook James O. McKinsey, Arley D. Cathey, Bloomberg L.P. CEO Daniel Doctoroff, Credit Suisse CEO Brady Dougan, Morningstar, Inc. founder and CEO Joe Mansueto, Chicago Cubs owner and chairman Thomas S. Ricketts, and NBA commissioner Adam Silver.*

Perturbed Paragraph: *In business, notable alumni include Microsoft CEO Satya Nadella, Oracle Corporation founder and the third richest man in America Larry Ellison, Goldman Sachs and MF Global CEO as well as former Governor of New Jersey Jon Corzine, McKinsey & Company founder and author of the first management accounting textbook James O. McKinsey, co-founder of the Blackstone Group Peter G. Peterson, co-founder of AQR Capital Management Cliff Asness, founder of Dimensional Fund Advisors David Booth, founder of The Carlyle Group David Rubenstein, Lazard CEO Ken Jacobs, entrepreneur David O. Sacks, CEO of TPG Group and former COO of Goldman Sachs Jon Winkelreid, former COO of Goldman Sachs Andrew Alper, billionaire investor and founder of Oaktree Capital Management Howard Marks, Bloomberg L.P. CEO Daniel Doctoroff, Credit Suisse CEO Brady Dougan, Morningstar, Inc. founder and CEO Joe Mansueto, Chicago Cubs owner and chairman Thomas S. Ricketts, and NBA commissioner Adam Silver.*

Question: What Goldman Sachs CEO is also an alumni of the University of Chicago?

Prediction of GPT-3.5-turbo-0125 and Llama-3-8B-Instruct: Jon Corzine→Jon Winkelreid

Prediction of Falcon-40B-Instruct: Jon Corzine→David Rubenstein, co-founder of The Carlyle Group, is also an alumnus of the University of Chicago.

NAT_V2_CHALLENGE

Original Paragraph: *Each chapter has a number of authors who are responsible for writing and editing the material. A chapter typically has two "coordinating lead authors", ten to fifteen "lead authors", and a somewhat larger number of "contributing authors". The coordinating lead authors are responsible for assembling the contributions of the other authors, ensuring that they meet stylistic and formatting requirements, and reporting to the Working Group chairs. Lead authors are responsible for writing sections of chapters. Contributing authors prepare text, graphs or data for inclusion by the lead authors.*

Perturbed Paragraph: *Each chapter has a number of authors to write and edit the material. A typical chapter has two coordinating lead authors, ten to fifteen lead authors and a larger number of contributing authors. The coordinating lead authors assemble the contributions of the other authors. They ensure that contributions meet stylistic and formatting requirements. They report to the Working Group co-chairs. Lead authors write sections of chapters. They invite contributing authors to prepare text, graphs or data for inclusion.*

Question: Who has the responsibility for publishing materials?

Prediction of Mistral-7B-Instruct-v0.2: Unanswerable. The text does not mention any responsibility related to publishing materials.→The coordinating lead authors are responsible for publishing materials in the given context.

Figure 4: Natural perturbed MRC examples that confuse LLMs.

Dataset	Victim	Attacker	CharOCR	CharInst	CharSubst	CharSwidMid	CharSwidRand	Winsert (CWE)	Wsubstite (CWE)	Wsplit	Wswap	Wdelete	Wcrop	WbySubst	Whitert (WE)	AddBert	DecPam	Style Transfer		
SQuAD	d141bert-base	bert-base-cased	8.71	8.27	8.18	8.28	8.28	3.25	3.05	3.17	1.33	0.78	1.46	14.60	14.60	7.78	7.38			
	bert-base-cased	bert-base-cased	-6.65	-6.82	-6.04	-5.95	-5.95	-2.17	-2.05	-2.29	-1.32	-0.67	-1.82	-1.32	-22.95	-22.95	-7.0	-6.65		
	bert-large-cased	bert-large-cased	-7.85	-7.32	-7.3	-6.91	-6.91	-2.18	-2.18	-2.18	-2.2	-1.44	-1.54	-1.44	-23.99	-24.08	-6.15	-6.55		
	bert-large-cased	bert-large-cased	-4.34	-4.06	-4.06	-4.2	-4.2	-2.13	-2.13	-2.13	-1.88	-1.02	-1.38	-1.02	-12.67	-12.67	-6.7	-5.55		
	spartan-large-cased	spartan-large-cased	-5.78	-5.31	-5.31	-5.06	-5.06	-2.21	-2.21	-2.21	-2.38	-1.51	-1.84	-1.51	-18.55	-18.55	-6.2	-5.91		
	roberta-base	roberta-base	-5.42	-5.66	-5.54	-5.49	-5.49	-2.14	-2.14	-2.14	-2.37	-1.65	-1.88	-1.65	-18.32	-18.32	-6.41	-5.64		
	roberta-large	roberta-large	-3.77	-4.19	-4.19	-3.99	-3.99	-1.88	-1.88	-1.88	-1.84	-1.27	-1.44	-1.27	-12.18	-12.18	-6.45	-5.42		
	roberta-base	roberta-base	-2.24	-2.5	-2.5	-2.83	-2.83	-1.8	-1.8	-1.8	-1.44	-0.93	-1.08	-0.93	-14.19	-14.19	-5.65	-4.75		
	roberta-base	roberta-base	-4.72	-4.72	-4.72	-4.72	-4.72	-2.1	-2.1	-2.1	-2.1	-1.09	-1.19	-1.09	-10.91	-10.91	-5.65	-4.75		
	roberta-base-v2	roberta-base-v2	-5.65	-6.47	-6.47	-6.27	-6.27	-2.1	-2.1	-2.1	-2.17	-1.89	-2.33	-1.89	-23.59	-23.59	-7.19	-6.18		
	roberta-base-v2	roberta-base-v2	-5.89	-6.27	-6.27	-6.03	-6.03	-2.08	-2.08	-2.08	-2.23	-1.54	-1.87	-1.54	-19.52	-19.52	-7.37	-6.18		
	albert-large-v1	albert-large-v1	-5.81	-5.81	-5.81	-6.05	-6.05	-1.94	-1.94	-1.94	-2.8	-1.29	-1.84	-1.29	-18.76	-18.76	-6.96	-6.14		
	albert-xl-large-v1	albert-xl-large-v1	-2.06	-2.44	-2.44	-2.82	-2.82	-1.7	-1.7	-1.7	-1.84	-1.12	-1.41	-1.12	-15.31	-15.31	-5.59	-4.68		
	albert-xxl-large-v1	albert-xxl-large-v1	-2.82	-3.18	-3.18	-3.56	-3.56	-1.6	-1.6	-1.6	-1.69	-1.06	-1.31	-1.06	-11.52	-11.52	-4.98	-4.45		
	deberta-large-v2	deberta-large-v2	-1.54	-1.78	-1.78	-2.24	-2.24	-1.36	-1.36	-1.36	-1.09	-0.7	-0.99	-0.7	-8.83	-8.83	-5.1	-4.57		
flan-t5-small	flan-t5-small	-7.93	-7.34	-7.34	-7.49	-7.49	-2.16	-2.16	-2.16	-2.88	-1.49	-1.61	-1.49	-10.02	-10.02	-7.85	-7.96			
flan-t5-base	flan-t5-base	-3.32	-4.36	-4.36	-4.8	-4.8	-2.06	-2.06	-2.06	-2.7	-1.71	-1.61	-1.71	-9.67	-9.67	-7.05	-6.03			
flan-t5-large	flan-t5-large	-1.24	-1.5	-1.5	-1.82	-1.82	-1.08	-1.08	-1.08	-1.01	-0.69	-0.85	-0.69	-5.97	-5.97	-5.18	-4.57			
flan-t5-xxl	flan-t5-xxl	-0.73	-0.98	-0.98	-1.24	-1.24	-0.89	-0.89	-0.89	-1.0	-0.7	-0.85	-0.7	-5.7	-5.7	-5.7	-3.82			
GPT-3.5-turbo-0125	GPT-3.5-turbo-0125	-0.34	-0.59	-0.59	-0.69	-0.69	-0.49	-0.49	-0.49	-0.7	-0.41	-0.57	-0.41	-10.2	-10.2	-5.7	-5.7			
Gemma-7B-IT	Gemma-7B-IT	-5.46	-6.53	-6.53	-7.9	-7.9	-3.62	-3.62	-3.62	-4.74	-4.53	-3.87	-4.53	-39.87	-39.87	-11.43	-10.29			
Gemma-7B-IT	Gemma-7B-IT	-2.17	-2.59	-2.59	-3.06	-3.06	-1.38	-1.38	-1.38	-1.76	-1.04	-1.31	-1.04	-14.97	-14.97	-11.6	-10.29			
Llama-2-7B	Llama-2-7B	-1.46	-1.79	-1.79	-2.17	-2.17	-1.46	-1.46	-1.46	-1.83	-1.19	-1.41	-1.19	-10.79	-10.79	-5.4	-4.79			
Llama-2-7B-Chat	Llama-2-7B-Chat	-0.76	-1.17	-1.17	-1.55	-1.55	-0.93	-0.93	-0.93	-1.54	-0.94	-1.19	-0.94	-6.97	-6.97	-5.1	-4.37			
Llama-3-8B-Inst	Llama-3-8B-Inst	-0.49	-0.73	-0.73	-1.13	-1.13	-0.34	-0.34	-0.34	-0.99	-0.64	-0.83	-0.64	-4.42	-4.42	-4.42	-3.83			
Mistral-7B-Inst	Mistral-7B-Inst	-0.42	-0.66	-0.66	-0.96	-0.96	-0.31	-0.31	-0.31	-0.86	-0.51	-0.77	-0.51	-3.87	-3.87	-4.42	-3.83			
Falcon-40B-Inst	Falcon-40B-Inst	-0.87	-1.36	-1.36	-1.72	-1.72	-1.1	-1.1	-1.1	-1.77	-1.13	-1.41	-1.13	-9.95	-9.95	-4.89	-2.75			
Falcon-40B-Inst	Falcon-40B-Inst	-1.83	-2.37	-2.37	-2.98	-2.98	-2.08	-2.08	-2.08	-2.81	-2.41	-3.07	-2.41	-12.98	-12.98	-8.01	-4.89			
SQuAD 2.0	distilbert-base	bert-base-cased	-2.65	-2.82	-2.82	-3.06	-3.06	-1.44	-1.44	-1.64	-0.64	-0.78	-0.64	-1.36	-1.36	-24.02	-24.02	-5.58	-4.28	
	bert-base-cased	bert-base-cased	-4.14	-4.64	-4.64	-5.05	-5.05	-3.21	-3.21	-3.21	-3.21	-1.76	-2.41	-1.76	-33.02	-33.02	-6.4	-5.94		
	bert-large-cased	bert-large-cased	-7.85	-7.32	-7.32	-6.91	-6.91	-2.18	-2.18	-2.18	-2.2	-1.44	-1.54	-1.44	-23.99	-24.08	-6.15	-6.55		
	bert-large-cased	bert-large-cased	-4.34	-4.06	-4.06	-4.2	-4.2	-2.13	-2.13	-2.13	-1.88	-1.02	-1.38	-1.02	-12.67	-12.67	-6.7	-5.55		
	spartan-large-cased	spartan-large-cased	-5.78	-5.31	-5.31	-5.06	-5.06	-2.21	-2.21	-2.21	-2.38	-1.51	-1.84	-1.51	-18.55	-18.55	-6.2	-5.91		
	roberta-base	roberta-base	-5.42	-5.66	-5.54	-5.49	-5.49	-2.14	-2.14	-2.14	-2.37	-1.65	-1.88	-1.65	-18.32	-18.32	-6.41	-5.64		
	roberta-large	roberta-large	-3.77	-4.19	-4.19	-3.99	-3.99	-1.88	-1.88	-1.88	-1.84	-1.27	-1.44	-1.27	-12.18	-12.18	-6.45	-5.42		
	roberta-base	roberta-base	-2.24	-2.5	-2.5	-2.83	-2.83	-1.8	-1.8	-1.8	-1.44	-0.93	-1.08	-0.93	-14.19	-14.19	-5.65	-4.75		
	roberta-base	roberta-base	-4.72	-4.72	-4.72	-4.72	-4.72	-2.1	-2.1	-2.1	-2.1	-1.09	-1.19	-1.09	-10.91	-10.91	-5.65	-4.75		
	roberta-base-v2	roberta-base-v2	-5.65	-6.47	-6.47	-6.27	-6.27	-2.1	-2.1	-2.1	-2.17	-1.89	-2.33	-1.89	-23.59	-23.59	-7.19	-6.18		
	roberta-base-v2	roberta-base-v2	-5.89	-6.27	-6.27	-6.03	-6.03	-2.08	-2.08	-2.08	-2.23	-1.54	-1.87	-1.54	-19.52	-19.52	-7.37	-6.18		
	albert-large-v1	albert-large-v1	-5.81	-5.81	-5.81	-6.05	-6.05	-1.94	-1.94	-1.94	-2.8	-1.29	-1.84	-1.29	-18.76	-18.76	-6.96	-6.14		
	albert-xl-large-v1	albert-xl-large-v1	-2.06	-2.44	-2.44	-2.82	-2.82	-1.7	-1.7	-1.7	-1.84	-1.12	-1.41	-1.12	-15.31	-15.31	-5.59	-4.68		
	albert-xxl-large-v1	albert-xxl-large-v1	-2.82	-3.18	-3.18	-3.56	-3.56	-1.6	-1.6	-1.6	-1.69	-1.06	-1.31	-1.06	-11.52	-11.52	-4.98	-4.45		
	deberta-large-v2	deberta-large-v2	-1.54	-1.78	-1.78	-2.24	-2.24	-1.36	-1.36	-1.36	-1.09	-0.7	-0.99	-0.7	-8.83	-8.83	-5.1	-4.57		
flan-t5-small	flan-t5-small	-7.93	-7.34	-7.34	-7.49	-7.49	-2.16	-2.16	-2.16	-2.88	-1.49	-1.61	-1.49	-10.02	-10.02	-7.85	-7.96			
flan-t5-base	flan-t5-base	-3.32	-4.36	-4.36	-4.8	-4.8	-2.06	-2.06	-2.06	-2.7	-1.71	-1.61	-1.71	-9.67	-9.67	-7.05	-6.03			
flan-t5-large	flan-t5-large	-1.24	-1.5	-1.5	-1.82	-1.82	-1.08	-1.08	-1.08	-1.01	-0.69	-0.85	-0.69	-5.97	-5.97	-5.18	-4.57			
flan-t5-xxl	flan-t5-xxl	-0.73	-0.98	-0.98	-1.24	-1.24	-0.89	-0.89	-0.89	-1.0	-0.7	-0.85	-0.7	-5.7	-5.7	-5.7	-3.82			
GPT-3.5-turbo-0125	GPT-3.5-turbo-0125	-0.34	-0.59	-0.59	-0.69	-0.69	-0.49	-0.49	-0.49	-0.7	-0.41	-0.57	-0.41	-10.2	-10.2	-5.7	-5.7			
Gemma-7B-IT	Gemma-7B-IT	-5.46	-6.53	-6.53	-7.9	-7.9	-3.62	-3.62	-3.62	-4.74	-4.53	-3.87	-4.53	-39.87	-39.87	-11.43	-10.29			
Gemma-7B-IT	Gemma-7B-IT	-2.17	-2.59	-2.59	-3.06	-3.06	-1.38	-1.38	-1.38	-1.76	-1.04	-1.31	-1.04	-14.97	-14.97	-11.6	-10.29			
Llama-2-7B	Llama-2-7B	-1.46	-1.79	-1.79	-2.17	-2.17	-1.46	-1.46	-1.46	-1.83	-1.19	-1.41	-1.19	-10.79	-10.79	-5.4	-4.79			
Llama-2-7B-Chat	Llama-2-7B-Chat	-0.76	-1.17	-1.17	-1.55	-1.55	-0.93	-0.93	-0.93	-1.54	-0.94	-1.19	-0.94	-6.97	-6.97	-5.1	-4.37			
Llama-3-8B-Inst	Llama-3-8B-Inst	-0.49	-0.73	-0.73	-1.13	-1.13	-0.34	-0.34	-0.34	-0.99	-0.64	-0.83	-0.64	-4.42	-4.42	-4.42	-3.83			
Mistral-7B-Inst	Mistral-7B-Inst	-0.42	-0.66	-0.66	-0.96	-0.96	-0.31	-0.31	-0.31	-0.86	-0.51	-0.77	-0.51	-3.87	-3.87	-4.42	-3.83			
Falcon-40B-Inst	Falcon-40B-Inst	-0.87	-1.36	-1.36	-1.72	-1.72	-1.1	-1.1	-1.1	-1.77	-1.13	-1.41	-1.13	-9.95	-9.95	-4.89	-2.75			
Falcon-40B-Inst	Falcon-40B-Inst	-1.83	-2.37	-2.37	-2.98	-2.98	-2.08	-2.08	-2.08	-2.81	-2.41	-3.07	-2.41	-12.98	-12.98	-8.01	-4.89			

Table 10: Robustness of MRC models under synthetic perturbations. For encoder-only and decoder-only models, the results displayed are the relative F1 change (%), while for decoder-only models, we demonstrate the relative IM change (%). In SQuAD 2.0, the values shown in the parentheses represent the relative change for answerable and unanswerable questions, respectively.

<p>CharSwapRand</p> <p>Original Paragraph: [...] According to the general principle of proportionality the lawfulness of an action depends on whether it was appropriate and necessary to achieve the objectives legitimately pursued. When there is a choice between several appropriate measures the least onerous must be adopted, and any disadvantage caused must not be disproportionate to the aims pursued. The principle of proportionality is also recognised in Article 5 of the EC Treaty, stating that "any action by the Community shall not go beyond what is necessary to achieve the objectives of this Treaty".</p> <p>Perturbed Paragraph: [...] According to the general principle of proportionality the lawfulness of an action depends on whether it was appaoptirre and necessary to achieve the ibjvotcees legitimately psdeuur. When there is a choice between several ipprapretoa seumraes the least onerous must be adopted, and any disadvantage caused must not be disproportionate to the aims upsured. The pricnplie of proportionality is also recognised in Article 5 of the EC Treaty, stating that " any acitno by the Community shall not go beyond what is necessary to achieve the objectives of this Treaty ".</p> <p>Question: Where is the principle of proportionality not recognized in the EC treaty?</p> <p>Prediction of GPT-3.5-turbo-0125: unanswerable→In Article 5 of the EC Treaty.</p> <p>Prediction by humans under perturbation: unanswerable</p>
<p>WSubstitute (CWE)</p> <p>Original Paragraph: One of the most famous people born in Warsaw was Maria SkÇodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the Nobel Prize. Famous musicians include WÇadysÇaw Szpilman and Fr@d©ric Chopin. Though Chopin was born in the village of elazowa Wola, about 60 km (37 mi) from Warsaw, he moved to the city with his family when he was seven months old. Casimir Pulaski, a Polish general and hero of the American Revolutionary War, was born here in 1745.</p> <p>Perturbed Paragraph: One of the most famous Russians born outside Warsaw was Maria SkÇodowska-Curie, who achieved international recognition for her research on radioactivity who was its last female recipient granted the Mercury Prize. Famous musicians were WÇadysÇaw Szpilman and Fr@d©ric Chopin. Though Chopin was born in the village of elazowa Wola, about 60 km (37 mi) from Warsaw, he moved to the city with his family when he only seven months old. Casimir Pulaski, a Polish general and hero of the American Revolutionary War, was married here in 1745.</p> <p>Question: Who was Fredderic Wola?</p> <p>Prediction of GPT-3.5-turbo-0125: unanswerable→Frederic Chopin.</p> <p>Prediction by humans under perturbation: unanswerable</p>
<p>WSplit</p> <p>Original Paragraph: Civil disobedience is usually defined as pertaining to a citizen's relation to the state and its laws, as distinguished from a constitutional impasse in which two public agencies, especially two equally sovereign branches of government, conflict. For instance, if the head of government of a country were to refuse to enforce a decision of that country's highest court, it would not be civil disobedience, since the head of government would be acting in her or his capacity as public official rather than private citizen.</p> <p>Perturbed Paragraph: C ivil disobedience is usually de fined as pertai ning to a citizen ' s relation to the st ate and its la ws, as distinguished from a constitutional impasse in which two pub lic age ncies, especially two eq ually sove reign b ranches of government, conflict. For ins tance, if the head of government of a c ountry were to refuse to en force a decision of that cou ntry ' s highest court, it would not be civil disobedience, since the he ad of government wo uld be acting in her or his capaci ty as p ublic offic ial rather t han private citizen.</p> <p>Question: What group promotes using conflict to decide cases?</p> <p>Prediction of GPT-3.5-turbo-0125: unanswerable→public agencies</p> <p>Prediction by humans under perturbation: unanswerable</p>

Figure 5: Valid synthetic adversarial examples (character and word levels).

<p>AddSent</p> <p>Original Paragraph: [...] <i>The Anglo-Norman language was eventually absorbed into the Anglo-Saxon language of their subjects (see Old English) and influenced it, helping (along with the Norse language of the earlier Anglo-Norse settlers and the Latin used by the church) in the development of Middle English. It in turn evolved into Modern English.</i></p> <p>Perturbed Passage: <i>What was the final outcome of the Anglo-Norman language’s influence on English literature? [...] The Anglo-Norman language was eventually absorbed into the Anglo-Saxon language of their subjects (see Old English) and influenced it, helping (along with the Norse language of the earlier Anglo-Norse settlers and the Latin used by the church) in the development of Middle English. It in turn evolved into Modern English.</i></p> <p>Question: What was the Anglo-Norman language’s final form?</p> <p>Prediction of GPT-3.5-turbo-0125: Modern English→It was eventually absorbed into the Anglo-Saxon language of their subjects.</p> <p>Prediction by humans under perturbation: Modern English</p>
<p>DocPara</p> <p>Original Paragraph: <i>The area is also known for its early twentieth century homes, many of which have been restored in recent decades. The area includes many California Bungalow and American Craftsman style homes, Spanish Colonial Revival Style architecture, Mediterranean Revival Style architecture, Mission Revival Style architecture, and many Storybook houses designed by Fresno architects, Hilliard, Taylor & Wheeler. The residential architecture of the Tower District contrasts with the newer areas of tract homes urban sprawl in north and east areas of Fresno.</i></p> <p>Perturbed Paragraph: <i>In recent decades, many early twentieth century homes in the north and east areas of Fresno have been restored, including California Bungalow and American Craftsman style properties, as well as Storybook houses designed by Hilliard, Taylor & Wheeler. The unique architectural styles of the Tower District, such as Spanish Colonial Revival, Mediterranean Revival, and Mission Revival, stand in contrast to the newer tract homes and urban sprawl in the surrounding areas.</i></p> <p>Question: Are California Bungalows located in the north or east?</p> <p>Prediction of GPT-3.5-turbo-0125: unanswerable→In the north and east areas of Fresno.</p> <p>Prediction by humans under perturbation: unanswerable</p>
<p>Style Transfer</p> <p>Original Paragraph: <i>Western musical instruments were introduced to enrich Chinese performing arts. From this period dates the conversion to Islam, by Muslims of Central Asia, of growing numbers of Chinese in the northwest and southwest. Nestorianism and Roman Catholicism also enjoyed a period of toleration. Buddhism (especially Tibetan Buddhism) flourished, although Taoism endured certain persecutions in favor of Buddhism from the Yuan government. Confucian governmental practices and examinations based on the Classics, which had fallen into disuse in north China during the period of disunity, were reinstated by the Yuan court, probably in the hope of maintaining order over Han society. Advances were realized in the fields of travel literature, cartography, geography, and scientific education.</i></p> <p>Perturbed Paragraph: <i>During the Yuan dynasty, Western musical instruments were introduced to enhance Chinese performing arts, while the influence of Islam from Central Asia led to the conversion of some Chinese individuals in the northwest and southwest regions. Nestorianism and Roman Catholicism were also tolerated during this period. Despite the flourishing of Buddhism, Taoism faced persecution by the Yuan government in favor of Buddhism. The revival of Confucian governmental practices and examinations based on the Classics aimed to maintain order in Han society. Additionally, significant progress was made in travel literature, cartography, geography, and scientific education during this time.</i></p> <p>Question: What fields of study were not advanced during the Yuan?</p> <p>Prediction of GPT-3.5-turbo-0125: unanswerable→Taoism</p> <p>Prediction by humans under perturbation: unanswerable</p>

Figure 6: Valid synthetic adversarial examples (sentence and document levels).

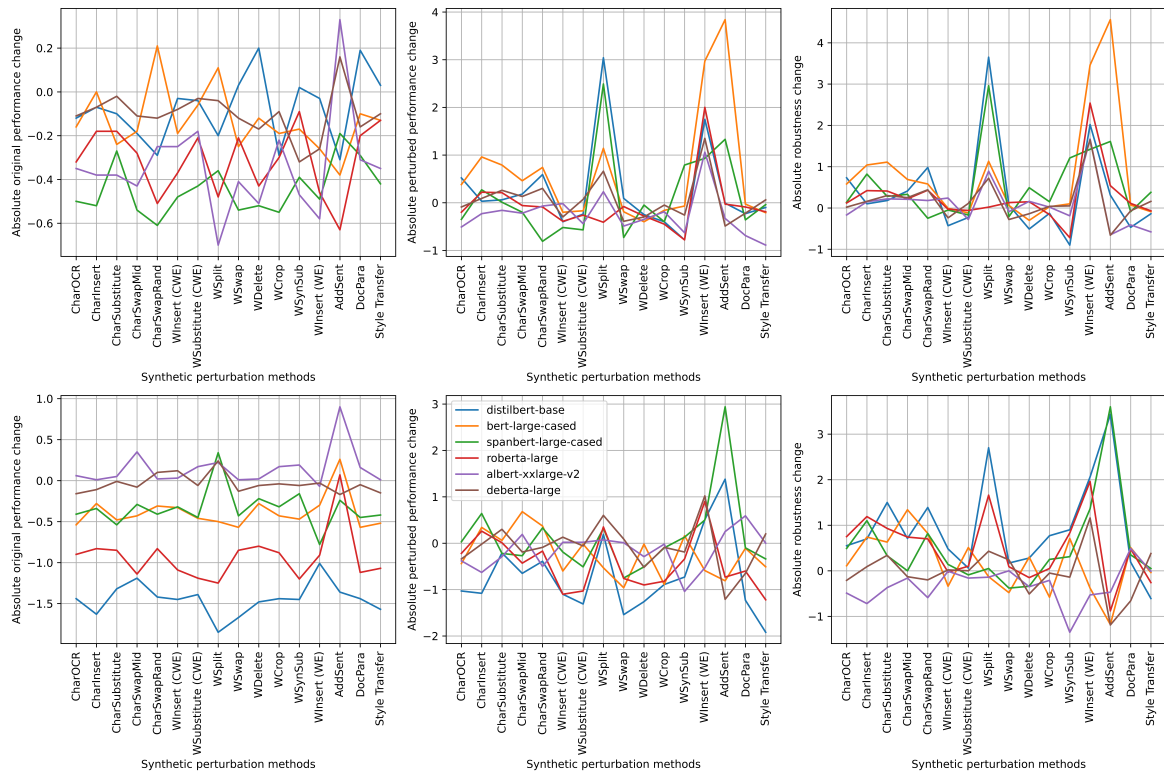


Figure 7: Absolute changes in original and perturbed performance (F1), as well as the robustness of six encoder-only models under various synthetic perturbations, following training on the augmented dataset with naturally perturbed MRC samples. The upper row of figures illustrates the outcomes obtained on the test sets created with SQuAD 1.1 as the reference, while the bottom row displays the results on SQuAD 2.0 format test datasets.