

Exploiting Positional Bias for Query-Agnostic Generative Content in Search

Andrew Parry
University of Glasgow

Sean MacAvaney
University of Glasgow

Debasis Ganguly
University of Glasgow

Abstract

In recent years, neural ranking models (NRMs) have been shown to substantially outperform their lexical counterparts in text retrieval. In traditional search pipelines, a combination of features leads to well-defined behaviour. However, as neural approaches become increasingly prevalent as the final scoring component of engines or as standalone systems, their robustness to malicious text and, more generally, semantic perturbation needs to be better understood. We posit that the transformer attention mechanism can induce exploitable defects through positional bias in search models, leading to an attack that could generalise beyond a single query or topic. We demonstrate such defects by showing that non-relevant text—such as promotional content—can be easily injected into a document without adversely affecting its position in search results. Unlike previous gradient-based attacks, we demonstrate these biases in a query-agnostic fashion. In doing so, without the knowledge of topicality, we can still reduce the negative effects of non-relevant content injection by controlling injection position. Our experiments are conducted with simulated on-topic promotional text automatically generated by prompting LLMs with topical context from target documents. We find that contextualisation of a non-relevant text further reduces negative effects whilst likely circumventing existing content filtering mechanisms. In contrast, lexical models are found to be more resilient to such content injection attacks. We then investigate a simple yet effective compensation for the weaknesses of the NRMs in search, validating our hypotheses regarding transformer bias.

1 Introduction

Neural Ranking Models (NRMs) have improved over lexical models on many retrieval benchmarks (Karpukhin et al., 2020; Nogueira et al., 2020). In contrast to lexical retrieval models such as TF-IDF weighting (Robertson and Jones, 1976)

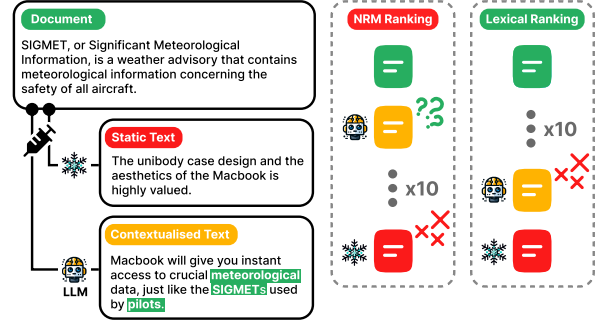


Figure 1: Injection of static and contextualised (document-conditioned by Llama 2) text into a document. BM25, aligned with retrieval axioms, penalises both injections; document length has increased with insufficient additional relevant information. However, the NRM (monoT5) is invariant to the addition of contextualised promotional text. Further examples are enlisted in Appendix E.

and BM25 (Robertson et al., 1998), pre-trained language models (PLMs) are pre-trained in an unsupervised manner over large corpora and fine-tuned on labelled examples, allowing for the encoding of text into a richly contextualised latent representation (Devlin et al., 2019; Liu et al., 2019) generally based on the transformer architecture (Vaswani et al., 2017). Besides the traditional “ad-hoc” setting, NRMs are applied in various NLP tasks, such as retrieval-augmented-generation (Lewis et al., 2020) in which a retrieval pipeline provides information to a downstream system to improve factuality. Where previously a search pipeline composed of multiple features, including lexical models, would be robust to the weaknesses of neural models such as adversarial attacks (Goodfellow et al., 2015; Wu et al., 2022) or weak term matching (Formal et al., 2022), we consider the robustness of NRMs as standalone retrieval models. If the relevance of a text varies simply by changing the ordering of constituent spans or the addition of non-relevant content is ignored by a model (MacAvaney et al., 2022), this behaviour diverges from

axioms outlining trustworthy and expected model response (Fang et al., 2011) and has potential to degrade user experience or propagate harmful information in an otherwise helpful text. Considering positional bias, prior work has noted the sensitivity of NRMs to span order (Jiang et al., 2021; MacAvaney et al., 2022); we, for the first time, consider how this bias may be exploited by malicious actors to understand better why NRMs exhibit such behaviour and how this limitation can be addressed.

To simulate a harmful adversary exploiting both positional and contextual bias in content injection, we apply large language models (LLMs) to condition the generation of a non-relevant span on a targeted document. Consider that an individual could promote, for example, a product or a political idea within the context of many different topics without requiring human intervention using generative models as shown in Figure 1. We condition our attack solely on document text, giving potential to an attack that is query agnostic. We propose a framework for the automated generation and injection of non-relevant content to maintain the rank of a document whilst completing an ulterior task.

We posit that **a)** the position of an injected span can have largely disproportionate effects on relevance when no change in information need has occurred, and **b)** by considering the context surrounding a span, we can better ‘hide’ this text and its effect on relevance estimation. We then empirically investigate the injection of text spans with variable relevance into passages, finding that injection position can largely affect the rank of documents. We observe that the rank of a document augmented with non-relevant text can be improved by up to 9 ranks simply by changing where the span is placed. From this observation, we propose the concept of *attention bleed-through*, the propagation of some positive context from one sequence to another. We then apply our findings to the ulterior task of promotion, finding that by exploiting the components of the transformer, we can reduce the effects of non-relevant text injection, improving over real examples of promotional content out of context, which we call static text as shown in Figure 1. We then propose an approach to mitigate the effects of document-conditioned generation motivated by our investigation of the above hypotheses. We recover significant ranking performance across multiple retrieval architectures without access to generated content.

In summary, NRMs present relevance-dependent positional bias, which, combined with document-conditioned non-relevant text, could not only degrade user experience but also propagate harmful content. As lexical models suitably penalise such content, our hypotheses and findings present wider implications for the application of NRMs outside of a sanitised evaluation setting. In the interest of reproducibility, we provide our data, artefacts, and scripts¹.

2 Preliminaries and Related Work

Neural ranking models often use fine-tuned pre-trained language models (Karpukhin et al., 2020; Nogueira et al., 2020) with strong natural language understanding. Both query–document representation (Nogueira et al., 2020; Pradeep et al., 2022) and separate representation (Karpukhin et al., 2020; Khattab and Zaharia, 2020) approaches exist. In BERT-based models, the representation of the [CLS] token is commonly used for relevance estimation between a query and a document with prior work finding positional bias when applying [CLS]-pooling (Jiang et al., 2021). However, multiple approaches to contextualised representations within retrieval exist beyond [CLS]-pooling, which we posit can also be affected by positional bias. Crucially, we target these contextualised representations as part of our investigation.

To exploit context in semantic search, the primary component of our workflow is a generative Large Language Model (LLM) (Radford et al., 2018; Touvron et al., 2023). Given an input sequence of words, an LLM outputs a weighted vector over a vocabulary of tokens, conditioning each subsequent generation on prior realizations (Radford et al., 2018). Recent works have applied instruction fine-tuning to LLMs, observing success in conversational question answering, among other NLP tasks. The description of a task to an LLM is commonly called prompt learning (Brown et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023). These natural language descriptions allow for the design of complex tasks and have shown strong zero-shot performance across a range of tasks that previously required a trained model (Shen et al., 2023; Li et al., 2023; Kojima et al., 2023, *inter alia*). It is this generalisation that we look to exploit to automatically produce non-relevant content in the context of a document.

¹ [Repository Link](#)

Prior works have probed the explainability of NRMs and behaviour that diverges from the desired retrieval objective or a human notion of relevance. [Câmara and Hauff \(2020\)](#) investigated how BERT-based ranking models’ relevance approximations correlate with retrieval axioms acting as the foundations for sparse term weighting models, finding that NRM relevance scores do not align with axiomatic approaches. [MacAvaney et al. \(2022\)](#) proposed the ABNIRML score, which measures NRM preference for one document set versus an augmented set. Notably, they found unexpected behaviours, such as NRM preference for a document with non-relevant text appended over the original document. Further studies show that NRMs fail to detect important lexical matches ([Formal et al., 2022](#)) and are invariant to semantics such as negation ([Weller et al., 2023](#)).

A common approach to exploit neural models is the adversarial attack. In this setting, a model input is perturbed with respect to the gradient response of either a target model ([Szegedy et al., 2014](#); [Goodfellow et al., 2015](#)) or a distilled surrogate ([Lord et al., 2022](#)). Prior works have considered attacking NRMs as a form of search engine optimisation, looking to improve the rank of a relevant document by modifying said document with respect to model output. Synonyms of words within the targeted document are substituted in an iterative process to approximate an optimal perturbation ([Wu et al., 2022](#); [Liu et al., 2023](#)). In the spirit of these attacks, our generation process ensures that the text is likely to be semantically coherent within the document’s body. However, reduce the effect of text injection in a query-agnostic fashion whilst achieving an ulterior objective as opposed to the sole objective of improving the rank of known relevant documents.

3 Methodology

This section outlines our methodology to generate contextual text for probing neural retrieval models.

Positional Information Transformer-based models are composed of a sequence of blocks containing the attention mechanism ([Vaswani et al., 2017](#)), which progressively manipulate the representation of a sequence. A positional encoding either informs the model of a token’s absolute position in a sequence ([Vaswani et al., 2017](#)) or its relative position compared to other tokens ([Shaw et al., 2018](#)). Though important for tasks such as text generation and, more generally, language modelling ([Shaw](#)

[et al., 2018](#)), we posit that this contextualisation conditioned on positional information can introduce unintended consequences for retrieval.

Given a sequence of tokens estimated to be relevant, we consider that through positional information, the attention mechanism which allows contextualisation between tokens could propagate “positive” attention scores from said relevant span to a subsequent non-relevant span. Formally, we hypothesise that for a sequence of two texts $\{t_0, t_1\}$ and a query q , if t_1 is considered more relevant than t_0 :

$$\mathcal{S}(q, t_1 \oplus t_0; \theta) > \mathcal{S}(q, t_0 \oplus t_1; \theta) \quad (1)$$

where in Equation 1, relevance score of a text t with respect to q by model θ is defined as $\mathcal{S}(q, t; \theta)$ and \oplus represents string concatenation to a sequence. We call these effects *attention bleed-through*, as in the context of retrieval, the positive effect of a relevant sentence ‘bleeds’ onto other sentences, reducing the negative effect on the downstream relevance score even if a span lacks inherent value to a task. This creates an unintended correlation between a span’s position and perceived relevance, potentially harming retrieval effectiveness.

To investigate this hypothesis more concretely, we investigate how adding non-relevant text close to document sentences with high similarity scores with the query affects retrieval. More concretely, we introduce *salience* as the scoring of a span within a text to another, e.g., a document to a query. If we decompose a document d into n spans, $d = \{s_0, \dots, s_n\}$, the most *salient span* with respect to some query q is defined as $\max_{s \in d} \text{sim}(e_q, e_s)$ where $\text{sim}(e_q, e_s)$ is a similarity measure in vector space e.g inner product.

Contextualised Generation Retrieval axioms suggest that when the length of a document increases without further satisfying an information need, it should decrease the rank of that document for a given query (TFC2, LNC1) ([Fang et al., 2011](#)). Traditionally, under lexical approaches ([Robertson et al., 1998](#)), adding non-relevant text decreases a document’s rank based on principles such as document length and information need mismatch. However, we propose that controlling the position and conditioning of promotional text can mitigate its negative impact.

We condition the generation of promotional spans on the context of the target document. This is

achieved using a generative language model, which receives a prompt specifying the promotional task and the document content itself. We hypothesise that this context will lead to partially shared topicality as the generated promotional text will likely contain tokens similar to those already present in the document, fostering positive attention interactions with the surrounding content. This contextual similarity may effectively obfuscate the promotional text within the document, minimising the ranking penalty imposed by retrieval models due to its non-relevance.

For a ranking model θ , we define the set of documents retrieved for a query q by the model as $\mathcal{R}(q; \theta)$. We condition the generation of each span on a target document such that relevance judgements are not required, creating a query-agnostic attack. Our prompt contains the task text sequence t , which, in this case, states that the model should promote an entity provided as input.

The output by a generative model ω of a document-conditioned span s within document $d \in \mathcal{R}(q; \theta)$, is governed by the conditional probability

$$P_\omega(s|t, d) \equiv \prod_{i=0}^{|s|} P_\omega(s_i|s_{0:i}, t, d) \quad (2)$$

where in Equation 2, $s_{0:i} = \{s_0, \dots, s_{i-1}\}$ are previously generated tokens. Tokens are realised by applying argmax over logits at each step i . Each token generation relies on the previously generated sequence, the task prompt, and the document context, aiming to create task-specific text that integrates within the document.

Mitigation within Retrieval Our experiments confirm that transformers are susceptible to manipulation through the position and contextualisation of injected text. This raises concerns about the general robustness of NRMs against similar attacks. Retraining each NRM may not be feasible. Instead, we propose a simple and scalable mitigation strategy based on our understanding of attention bleed-through. The core principle is that transformer-based models struggle to penalise non-relevant text due to potential bleed-through effects. We introduce a classification model ϕ specifically trained to detect promotional content within documents. This model operates independently of the ranking model.

We exploit our hypothesis of attention bleed-through to our advantage, considering that a slid-

ing window scoring method may minimise opportunities for bleed-through. By processing text in sliding windows, the classifier can focus on local interactions between potentially promotional content and its surrounding text. This isolation assists the classifier in making clearer estimations of the presence of undesirable text, potentially mitigating bleed-through’s negative impact. As the posterior represents a higher probability of promotion being present, we look to improve the relevance score if this estimate is lower. We subtract the posterior from 1 such that $\mathcal{S}'(d; \phi) = \max_{s \in d} 1 - P_\phi(Y = 1|s)$ reflects the posterior of the span s *not* containing promotion. As such, we determine the final relevance score of a document as follows:

$$\mathcal{S}(q, d) = \alpha \mathcal{S}(q, d; \theta) + (1 - \alpha) \mathcal{S}'(d; \phi) \quad (3)$$

In Equation 3, scoring combines the original NRM output parameterised by θ with the promotional content detection score parameterised by ϕ . The parameter α balances relevance and penalisation of undesirable content (based on promotional content detection). This approach offers a flexible trade-off between these competing objectives.

Research Questions We investigate the following research questions to validate our hypotheses.

- **RQ-1:** When injecting spans of text, how is relevance affected by position and context?
- **RQ-2:** Can we defend against contextualised text injection with a model-agnostic approach?

4 Evaluation of Contextualised Text

We now outline the evaluation setup of our approach and discuss findings from empirical evidence.

IR Datasets We use the MSMARCO passage collection (Bajaj et al., 2016), a corpus of over 8.8 million passages extracted from Bing query searches. In all experiments, we evaluate target models on the TREC Deep Learning 2019 track test set (Craswell et al., 2020), which contains relevance judgements for 43 queries providing human-assessed queries to validate our hypotheses. We use these human-judged relevance labels to inject spans of varying relevance.

Models We investigate NRMs based on BERT and T5 architectures to assess how both embedding- and decoder-determined relevance approximations

are affected by our approach. Additionally, we employ BM25 to contrast against neural models and to act as an indicator of exact term matching between queries and generated text. We evaluate two bi-encoders, Contriever (Izacard et al., 2022) and ColBERT (Khattab and Zaharia, 2020), which uses the late interaction paradigm. We also evaluate a sequence-to-sequence cross-encoder, monoT5 (Nogueira et al., 2020). A description of these models can be found in Appendix C.2.

Promotion Datasets In determining entities to promote, we chose five spans that explicitly referenced an entity from a subset of scraped Wikipedia rejected edits (Bertsch and Bethard, 2021). We consider these edits to be examples of static text as they represent promotional content outside the context of a document because human editors rejected these edits specifically for being promotional. The referenced entity becomes part of our task text in prompting each language model (as outlined in Section 3).

We inject promotional spans into the top 100 ranked documents for each query. We truncate each generation to a single sentence so that the span length does not confound comparisons. Each span is injected as described in Section 4.

Metrics We perform a pairwise evaluation comparing injected static text and document contextualised text, investigating how the negative effects of text injection on retrieval performance can be reduced from the perspective of a malicious content provider. We first evaluate the ABNIRML score (ABN) being the listwise preference of a model for two rankings. When the ABNIRML score is positive, the model prefers the original set; when it is 0, the augmentation does not affect preference, and a negative score indicates a preference for the augmented set. To assess the magnitude of the effects of text injection, we measure Mean Rank Shift (MRS) as the mean of the difference in rank between each document with and without injected text. A full description of these metrics is provided in Appendix C.1. In each evaluation, significance is determined by a 95% confidence t-test with Bonferroni correction to assess the significance of positional bias.

Injection of Known Relevant and Non-Relevant Text Investigating RQ-1, we evaluate ranking model preference for documents with and without injected text. Re-ranking is performed for each test

query over 100 documents retrieved by BM25. We inject spans from varying levels of human-judged relevant documents (ranging from 0-3) to assess the effect of on-topic and off-topic injection on relevance. Per common conventions with the DL19 dataset, we consider a judgement of 2-3 as *relevant*, 0-1 as *related*, and judgments from other queries as *non-relevant*. We inject documents at the following positions. We inject at different positions to assess how positional bias affects retrieval score when controlling for the known relevance of a document.

- **Absolute Position (Abs-P):** places a span of text before, in the middle or after a document
- **Relative Position (Rel-P):** places a span before and after the most salient sentence in a document (salience as defined in Section 3)

In particular, for obtaining embedded representations of sentences used for salience computation (defined in Section 3), we use the sentence transformers (Reimers and Gurevych, 2019) MPNET encoder model², which uses mean pooling over token embeddings.

Generation of Contextualised Promotional Text

We use a zero-shot approach to prompt a generative language model.

Prompt Format: Using the important keywords taken from the Document, write a sentence mentioning and promoting the Item:
Document: {*document*}
Item: {*item*}
Response:

We apply two LLMs of significantly varying sizes, Llama-2³ with 7 billion parameters and GPT-3.5 turbo⁴ with 175 billion parameters accessed through the OpenAI API. In investigating Llama 2, we show the feasibility of this approach inexpensively by running inference on a single GPU, which allows for the reproducibility of our findings, given that the underlying GPT-3.5 model frequently changes.

4.1 Results and Discussion

We now enlist the main observations from our experiments as follows.

² sentence-transformers/all-mpnet-base-v2

³ meta-llama/Llama-2-7b-chat-hf ⁴ gpt-3.5-turbo

Table 1: Wikipedia Entities and static examples of promotion. Examples are taken from Wikipedia edits rejected for being considered promotional.

Entity	Static Promotion
Finlandia Vodka	Drinkers view Finlandia vodka as a prestigious, reputable and great tasting brand.
Honda Motorcycles	Honda’s advance in western motorcycle markets of the 1960s was noted for its speed and power as well as its reliability.
Russia	Russia has a rich material culture and tradition in technology.
Macbook	The unibody case design and the aesthetics of the Macbook is highly valued.
Czech Republic	The Czech Republic was described by the guardian as one of Europe’s most flourishing economies.

Table 2: Injecting of spans from relevant, related, and non-relevant documents measuring ABNIRML (\downarrow) and MRC (\downarrow). Significance comparing positional injection Before and After.

		BM25		ColBERT		monoT5		Contriever	
		ABN	MRC	ABN	MRC	ABN	MRC	ABN	MRC
Absolute Position									
Before	Relevant	0.038	-5.161	-0.356	-16.760	-0.434	-18.136	-0.492	-17.107
	Related	0.199	-0.906	-0.099	-7.296	-0.106	-7.517	-0.230	-7.560
	N-Relevant	0.972	18.073	0.536	5.820	0.664	9.116	0.666	13.600
Middle	Relevant	0.038	-5.161	-0.462	-15.712	-0.541	-17.042	-0.533	-14.119
	Related	0.199	-0.906	-0.246	-7.797	-0.263	-8.187	-0.304	-7.111
	N-Relevant	0.972	18.073	0.207	1.542	0.493	4.666	0.494	7.544
After	Relevant	0.038	-5.161	-0.522 [†]	-14.525 [†]	-0.653 [†]	-15.860 [†]	-0.601 [†]	-12.394 [†]
	Related	0.199	-0.906	-0.312 [†]	-7.669	-0.431 [†]	-8.488	-0.402 [†]	-6.929
	N-Relevant	0.972	18.073	0.091 [†]	-0.348 [†]	0.297 [†]	1.912 [†]	0.360 [†]	4.863 [†]
Relative Position									
Before	Relevant	0.038	-5.161	-0.395	-15.536	-0.475	-16.917	-0.507	-14.588
	Related	0.199	-0.906	-0.160	-6.888	-0.183	-7.474	-0.253	-6.441
	N-Relevant	0.972	18.073	0.344	3.651	0.530	6.867	0.548	10.429
After	Relevant	0.038	-5.161	-0.529 [†]	-15.375	-0.607 [†]	-16.747	-0.586 [†]	-13.483
	Related	0.199	-0.906	-0.333 [†]	-8.240 [†]	-0.365 [†]	-8.778 [†]	-0.380 [†]	-7.358
	N-Relevant	0.972	18.073	0.070 [†]	-0.010 [†]	0.375 [†]	3.042 [†]	0.375 [†]	5.762 [†]

Position largely affects relevance under augmentation. In Table 2, we observe that across all neural models, position has a large effect on relevance (compare ‘before’ and ‘after’). BM25 penalises the addition of non-relevant content, reducing rank by 18 places as no query terms are likely present in the new span. Contriever shows minimal ABNIRML preference and rank change when appending non-relevant spans to a document. It is most likely that this results from the maximised distance between the [CLS] token, leading to reduced change in representation. Moreover, we observe similar invariance when probing ColBERT and monoT5, which do not pool their representations and use different positional encodings, suggesting that this effect can generalise across transformer variants. This observation correlates with our hypothesis outlined in Section 3, in that by injecting a span after the most salient text (‘after’, determined by both absolute and relative position), we reduce the effect of a sequence on the overall relevance approximation. Additionally, in Figure 2, observe

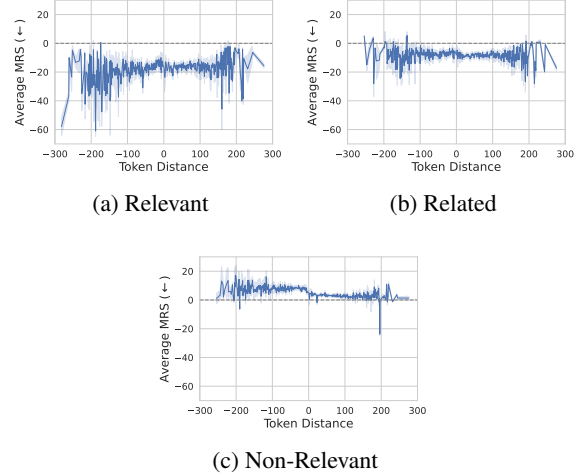


Figure 2: Average MRS by injected token distance to salient sequences for monoT5, observe noted reduction in variance near a salient sequence and clear reduction in penalty for non-relevant text after a salient sequence.

this variance in rank change reduces dramatically across all relevance levels when injecting near a salient span. This is notable as though absolute and relative positions show similar MRS, relative injection with respect to a salient span represents a more consistent reduction in rank penalty and, therefore, could be a more effective attack vector.

Across all NRMs, MRS is significantly reduced when injecting text directly after the most salient span compared to injection ‘before’. Though this position does not maximise absolute distance from the start of a document, the reduction in preference suggests that the influence of a sequence of tokens is partially determined by its position with respect to important tokens for a particular query. Furthermore, rank improvements by adding relevant spans are reduced when injected after the salient span. We observe that injection at the start of a document maximises the rank change, improving MRS by up to 18 ranks in the case of monoT5. This conforms to the findings of Liu et al. (2023) in that tokens that improve relevance (in their case, adversarial perturbations) are best placed at the start of a text. A small improvement in MRS is observed when appending non-relevant text across monoT5 and ColBERT, which is likely caused by a bias induced by over-fitting to the distribution of the corpus.

Conditioning generation on documents effectively reduces model preference and rank degradation. In Table 3, observe that BM25 rank and preference are reduced when injecting LLM-

generated text compared to a static promotional span with ABNIRML preference reducing from 0.99 to 0.58 in the case of GPT-generated promotion, and MRS is improved from 16.97 to 10.48. This occurs because of the topical contextualisation of the injected text as illustrated in Figure 1. We observe similar performance in Llama 2, with both models improving ABNIRML preference. ColBERT and monoT5 both show a strong bias for contextualised text, likely due to their stronger term matching. Though ABNIRML preference is reduced comparing static and contextualised text in Contriever, larger changes in MRS are mainly determined by the injection position with the model showing relative invariance to the context of an injection, suggesting that any exact term matching is minimal (Formal et al., 2022).

Effective positional injection amplifies contextualisation. In the case of both ColBERT and Contriever, when text is injected before either the document or the most salient span (Rows containing Before, columns ColBERT and Contriever), the models are invariant to the context of injected spans as the text has a large influence on the text which follows it. MonoT5 shows a clear preference for generated promotion even when placed before the documents (0.84 ABNIRML versus 0.60 and 0.57), suggesting that its joint representation is susceptible to positive contextualisation between query and document terms. When combining document-conditioned generation and positional injection, rank change between the original and augmented sets can be reduced to an MRS of 1.86 for ColBERT and 1.80 for monoT5. Notably, though Contriever is susceptible to positional injection contrasting before and after in all cases when injecting contextualised promotion, ABNIRML preference is still reduced.

5 Mitigation Evaluation

We now discuss the evaluation and findings of RQ-2. Multiple architectures are affected by the methods described in Section 3. We aim to mitigate the adversarial effects of injected text without re-training an underlying search system. Thus, we use an independent model that avoids re-training existing NRMs.

Datasets As ground truth does not exist for this task and capturing the distribution of LLM-generated text is dependent on the particular

Table 3: Injection of promotional spans showing ABNIRML (ABN.) (\downarrow) and MRS (\downarrow). Significance comparing positional injection Before and After.

		BM25		ColBERT		monoT5		Contriever	
		ABN	MRS	ABN	MRS	ABN	MRS	ABN	MRS
Absolute Position									
Before	Static	0.993	16.971	0.786	9.579	0.830	12.867	0.671	12.286
	Llama-2	0.619	12.254	0.645	9.864	0.600	9.129	0.516	9.964
	GPT-3.5	0.576	10.483	0.667	9.513	0.570	8.719	0.526	9.944
Middle	Static	0.993	16.971	0.528	5.529	0.646	6.855	0.454	6.348
	Llama-2	0.619	12.254	0.435	4.676	0.432	4.692	0.383	5.422
	GPT-3.5	0.576	10.483	0.461	4.603	0.431	4.610	0.391	5.452
After	Static	0.993	16.971	0.423 [†]	3.170 [†]	0.596 [†]	4.561 [†]	0.350 [†]	4.406 [†]
	Llama-2	0.619	12.254	0.318 [†]	1.694 [†]	0.246 [†]	1.589 [†]	0.267 [†]	2.780 [†]
	GPT-3.5	0.576	10.483	0.354 [†]	1.859 [†]	0.261 [†]	1.804 [†]	0.289 [†]	3.079 [†]
Relative Position									
Before	Static	0.993	16.971	0.610	7.469	0.679	9.345	0.513	8.747
	Llama-2	0.619	12.254	0.534	7.688	0.502	7.040	0.466	8.096
	GPT-3.5	0.576	10.483	0.570	7.587	0.493	6.828	0.468	8.017
After	Static	0.993	16.971	0.422 [†]	4.024 [†]	0.585 [†]	5.620 [†]	0.348 [†]	4.774 [†]
	Llama-2	0.619	12.254	0.341 [†]	2.818 [†]	0.306 [†]	2.884 [†]	0.287 [†]	3.698 [†]
	GPT-3.5	0.576	10.483	0.359 [†]	2.833 [†]	0.302 [†]	2.856 [†]	0.288 [†]	3.756 [†]

model (Mitchell et al., 2023; Su et al., 2023), we instead penalise promotion directly, though our approach could feasibly be extended to a broader set of undesirable texts. We assume no access to the output of a particular generative model. Hence, we use a classifier in a zero-shot manner. We employ the SemEval 2020 Task 11 Propaganda Techniques Corpus (PTC) (Da San Martino et al., 2020) comprising 15000 spans of text labelled with 18 propaganda classes. We observe that the 17 classes denoting propaganda are suitable for classifying promotion. We balance the training set after the combination of the 17 labels such that we have an equal distribution of propaganda and non-propaganda. The dataset contains challenging examples constituting ‘weasel’ words, i.e., subtle promotion methods or discrediting of an entity (Ott, 2018; Bertsch and Bethard, 2021), which we consider to be a useful parallel to promotion.

Training We finetune a RoBERTa base model (Liu et al., 2019) in a binary classification setting (described in Appendix C.3) to model the posterior probability of promotion being present by taking the classifier’s confidence. We use RoBERTa as prior work has succeeded in the main SemEval task leveraging RoBERTa (Raj et al., 2020; Singh et al., 2020).

To evaluate the intermediate task of identifying whether a given document contains promotional text, we create a balanced test set by sub-sampling examples from the different combinations of each generative model, injection position, and promoted entities. This intermediate evaluation can be found in Appendix A. We evaluate in a retrieval setting

Table 4: Evaluation of our mitigation strategy against contextualised text generated by GPT-3.5 (RQ-2) measuring nDCG@10 (\uparrow), MRR (\uparrow) and MRPR (\downarrow). Optimal α denoted α^* . Δ with respect to baseline retrieval. Significance with respect to retrieval performance with no mitigation is denoted with \dagger .

	α^*	nDCG@10 (Δ)	MRR (Δ)	MRPR (Δ)
BM25				
Abs-P Before	0.1	0.414 (+0.068) \dagger	0.619 (+0.014)	0.121 (-0.136) \dagger
Abs-P Middle	0.1	0.406 (+0.060) \dagger	0.616 (+0.011)	0.130 (-0.127) \dagger
Abs-P After	0.1	0.404 (+0.058) \dagger	0.616 (+0.011)	0.139 (-0.118) \dagger
Rel-P Before	0.1	0.409 (+0.063) \dagger	0.618 (+0.013)	0.130 (-0.127) \dagger
Rel-P After	0.1	0.405 (+0.059) \dagger	0.616 (+0.011)	0.134 (-0.123) \dagger
ColBERT				
Abs-P Before	0.3	0.648 (+0.029)	0.862 (+0.004)	0.119 (-0.063) \dagger
Abs-P Middle	0.1	0.614 (+0.060) \dagger	0.815 (-0.027)	0.112 (-0.167) \dagger
Abs-P After	0.1	0.607 (+0.097) \dagger	0.813 (-0.024)	0.126 (-0.208) \dagger
Rel-P Before	0.3	0.643 (+0.032)	0.862 (+0.004)	0.130 (-0.057) \dagger
Rel-P After	0.1	0.614 (+0.083) \dagger	0.816 (-0.014)	0.108 (-0.203) \dagger
monoT5				
Abs-P Before	0.9	0.634 (+0.019)	0.838 (-0.031)	0.143 (-0.072) \dagger
Abs-P Middle	0.6	0.611 (+0.051) \dagger	0.782 (-0.062)	0.127 (-0.183) \dagger
Abs-P After	0.6	0.606 (+0.100) \dagger	0.780 (-0.035)	0.134 (-0.272) \dagger
Rel-P Before	0.9	0.624 (+0.023)	0.838 (-0.031)	0.153 (-0.083) \dagger
Rel-P After	0.6	0.610 (+0.065) \dagger	0.782 (-0.053)	0.126 (-0.216) \dagger
Contriever				
Abs-P Before	0.2	0.589 (+0.054) \dagger	0.783 (+0.059)	0.151 (-0.147) \dagger
Abs-P Middle	0.3	0.573 (+0.079) \dagger	0.765 (+0.048)	0.159 (-0.207) \dagger
Abs-P After	0.4	0.567 (+0.096) \dagger	0.760 (+0.044)	0.139 (-0.242) \dagger
Rel-P Before	0.3	0.581 (+0.056) \dagger	0.768 (+0.049)	0.152 (-0.171) \dagger
Rel-P After	0.4	0.573 (+0.091) \dagger	0.761 (+0.043)	0.127 (-0.239) \dagger

described in Section 3 measuring nDCG@10 and MRR. We create a pseudo-corpus where one augmented document is added for each entity such that we have 200 documents for each query. Metrics are aggregated over each entity, simulating the scenario of an automated targeted promotion of each entity in search. Relevance judgements are provided for only augmented documents such that we can observe the penalisation of promotion by MRR, which we name **Mean Reciprocal Promotional Reduction (MRPR)** as we want to *minimise* this metric. We present the optimal value of α tuned on nDCG@10. Significance is determined as described in Section 4.

5.1 Results and Discussion

In Table 4, it can be seen that applying our mitigation recovers a significant fraction of retrieval performance in a standard setting. In 3 of 4 models, a large weighting of the classifier weight is required to successfully penalise promotion (generally $\alpha \sim [0.1, 0.3]$, Sensitivity analysis can be found in Appendix B). We observe cases of MRR being reduced by a statistically insignificant margin; given that the aggregate MRS of our augmentations degrades rank by 1 to 2 positions, one would

expect that MRR change is minimal from our previous experiments (Table 2). Model performance is still around 5 points of nDCG@10 lower than standard reported values; this is partially due to a fully duplicated corpus, leading to previously ‘highly-relevant’ documents which have been augmented with positional injection maintaining rank, as to provide a greater weight to the classifier would further reduce precision. Observe that the positional effects described in Section 4.1 are still present (see rows ‘after’); however, they are consistently reduced with nDCG@10 improved in the range [0.065, 0.101] for neural models and a significant reduction of MRPR occurs in all cases. In the case of monoT5, we see that a larger α is required to recover performance. Due to the binary nature of monoT5 relevance scoring (true and false), its relevance scoring is bi-modal; as such, higher values of α allow for the proper discrimination of promotion as all ‘relevant’ document scores are close. When considering a granular evaluation such as ABNIRML or MRS in Section 4.1, we see that absolute position is most effective; however, when evaluating known relevant documents, we see a similar negative effect on retrieval performance. This is notable in that for a document that meets an information need, and relative position is as effective under mitigation in maintaining rank as absolute position, further strengthening the discussion of results in Section 4.1.

6 Concluding Remarks

We have presented a novel investigation of text injection exploiting both relative positional injection and contextualisation via Large Language Models, reducing the negative effects of content injection on the rank of documents across multiple retrieval architectures. We propose the notion of *attention bleed-through* which we show to have implications for the robustness of NRMs. We then present model-agnostic mitigation, which improves nDCG@10 significantly under a classic evaluation setting by reducing the effect of contextualisation. We consider that these findings have wider implications in semantic search outside of a clean benchmark environment, arbitrarily changing relevance in a way that is not conducive to better aligning with information need.

Limitations

The clearest limitation of this approach is that we cannot yet *improve* the rank of a document with injected text. For example, this effect may be alleviated by promoting a particular entity within a context where it would be present naturally. However, this would provide far less insight into the weaknesses of NRMs as, most likely, a lexical model would also be affected. Furthermore, this approach depends on the generative model having some knowledge of the target entity. A tailored prompt that provides additional details of the entity could trivialise this point; however, for evaluation purposes, we have not tailored any prompt.

Ethics Statement

Our initial study has shown that what is considered a ‘small’ Large Language Model (7 billion parameters) can still contextualise within an abstract task in an effective way. As more LLM-generated content pollutes open text on platforms such as social media, we hypothesise that the automation of this process combined with prompts tuned to a particular entity or topic could pose problems for semantic search engines. We suggest that one cannot rely on generated text detection due to the many existing open models, so it could become infeasible to use model-specific checks.

We consider positional bias a core problem in neural IR. In a real situation, these strategies could be combined with more traditional adversarial methods to increase the rank of a document and minimise the undesirable effects of text, which achieves an ulterior objective, as shown in this work.

References

- Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2016. *MS MARCO: A Human Generated Machine Reading Comprehension Dataset*. *CEUR Workshop Proceedings*, 1773. Publisher: CEUR-WS.
- Amanda Bertsch and Steven Bethard. 2021. *Detection of Puffery on the English Wikipedia*. pages 329–333.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. *Language Models are Few-Shot Learners*. ArXiv:2005.14165 [cs].
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. *PaLM: Scaling Language Modeling with Pathways*. ArXiv:2204.02311 [cs].
- Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel Campos, and Ellen M. Voorhees. 2020. *Overview of the TREC 2019 deep learning track*.
- Arthur Câmara and Claudia Hauff. 2020. *Diagnosing BERT with Retrieval Heuristics*. In *Advances in Information Retrieval*, Lecture Notes in Computer Science, pages 605–618, Cham. Springer International Publishing.
- Giovanni Da San Martino, Alberto Barrón-Cedeño, Henning Wachsmuth, Rostislav Petrov, and Preslav Nakov. 2020. *SemEval-2020 Task 11: Detection of Propaganda Techniques in News Articles*. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 1377–1414, Barcelona (online). International Committee for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. ArXiv:1810.04805 [cs].
- Hui Fang, Tao Tao, and Chengxiang Zhai. 2011. *Diagnostic evaluation of information retrieval models*. *ACM Trans. Inf. Syst.*, 29(2).
- Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2022. *A study of lexical matching in neural information retrieval - abstract**. In *Proceedings of the 2nd Joint Conference of the Information*

- Retrieval Communities in Europe (CIRCLE 2022), Samatan, Gers, France, July 4-7, 2022*, volume 3178 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. [Explaining and Harnessing Adversarial Examples](#). ArXiv:1412.6572 [cs, stat].
- Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. [Efficiently Teaching an Effective Dense Retriever with Balanced Topic Aware Sampling](#). ArXiv:2104.06967 [cs].
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. [Unsupervised dense information retrieval with contrastive learning](#). *Trans. Mach. Learn. Res.*, 2022.
- Zhiying Jiang, Raphael Tang, Ji Xin, and Jimmy Lin. 2021. [How does BERT rerank passages? an attribution analysis with information bottlenecks](#). In *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 496–509, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [Dense Passage Retrieval for Open-Domain Question Answering](#). pages 6769–6781.
- Omar Khattab and Matei Zaharia. 2020. [ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT](#). *SIGIR 2020 - Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 39–48. ISBN: 9781450380164 Publisher: Association for Computing Machinery, Inc.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. [Large Language Models are Zero-Shot Reasoners](#). ArXiv:2205.11916 [cs].
- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Chen Li, Yixiao Ge, Jiayong Mao, Dian Li, and Ying Shan. 2023. [TagGPT: Large Language Models are Zero-shot Multimodal Taggers](#). ArXiv:2304.03022 [cs].
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [RoBERTa: A Robustly Optimized BERT Pretraining Approach](#). ArXiv:1907.11692 [cs].
- Yu-An Liu, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Wei Chen, Yixing Fan, and Xueqi Cheng. 2023. [Topic-oriented adversarial attacks against black-box neural ranking models](#).
- Nicholas A. Lord, Romain Mueller, and Luca Bertinetto. 2022. [Attacking deep networks with surrogate-based adversarial black-box methods is easy](#). ArXiv:2203.08725 [cs].
- Sean MacAvaney, Sergey Feldman, Nazli Goharian, Doug Downey, and Arman Cohan. 2022. [ABNIRML: Analyzing the Behavior of Neural IR Models](#). *Transactions of the Association for Computational Linguistics*, 10:224–239.
- Sean MacAvaney, Craig Macdonald, Charlie Clarke, Benjamin Piwowarski, and Harry Scells. IR Measures API. https://github.com/terrierteam/ir_measures/tree/main. Accessed: 2023-06-03.
- Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. [Detectgpt: Zero-shot machine-generated text detection using probability curvature](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 24950–24962. PMLR.
- Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. 2020. [Document Ranking with a Pre-trained Sequence-to-Sequence Model](#). pages 708–718.
- Douglas E. Ott. 2018. [Hedging, Weasel Words, and Truthiness in Scientific Writing](#). *JSLIS : Journal of the Society of Laparoendoscopic Surgeons*, 22(4):e2018.00063.
- Iadh Ounis, Gianni Amati, Vassilis Plachouras, Ben He, Craig Macdonald, and Douglas Johnson. 2005. Terrier information retrieval platform. In *Advances in Information Retrieval: 27th European Conference on IR Research, ECIR 2005, Santiago de Compostela, Spain, March 21-23, 2005. Proceedings 27*, pages 517–519. Springer.
- Ronak Pradeep, Yuqi Liu, Xinyu Zhang, Yilin Li, Andrew Yates, and Jimmy Lin. 2022. [Squeezing Water from a Stone: A Bag of Tricks for Further Improving Cross-Encoder Effectiveness for Reranking](#). In *Advances in Information Retrieval: 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10–14, 2022, Proceedings, Part I*, pages 655–670, Berlin, Heidelberg. Springer-Verlag.
- A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. 2018. [Improving language understanding with unsupervised learning](#).
- Mayank Raj, Ajay Jaiswal, Rohit R. R, Ankita Gupta, Sudeep Kumar Sahoo, Vertika Srivastava, and Yeon Hyang Kim. 2020. [Solomon at SemEval-2020 Task 11: Ensemble Architecture for Fine-Tuned Propaganda Detection in News Articles](#). ArXiv:2009.07473 [cs].

- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks](#). *EMNLP-IJCNLP 2019 - 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, Proceedings of the Conference*, pages 3982–3992. ISBN: 9781950737901 Publisher: Association for Computational Linguistics.
- Stephen E. Robertson and Karen Spärck Jones. 1976. [Relevance weighting of search terms](#). *J. Am. Soc. Inf. Sci.*, 27(3):129–146.
- Stephen E. Robertson, Steve Walker, and Micheline Hancock-Beaulieu. 1998. Okapi at TREC-7: automatic ad hoc, filtering, VLC and interactive. In *Proceedings of The Seventh Text REtrieval Conference, TREC 1998, Gaithersburg, Maryland, USA, November 9-11, 1998*, volume 500-242 of *NIST Special Publication*, pages 199–210. National Institute of Standards and Technology (NIST).
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. [Self-attention with relative position representations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 464–468, New Orleans, Louisiana. Association for Computational Linguistics.
- Tao Shen, Guodong Long, Xiubo Geng, Chongyang Tao, Tianyi Zhou, and Daxin Jiang. 2023. [Large Language Models are Strong Zero-Shot Retriever](#). ArXiv:2304.14233 [cs].
- Paramansh Singh, Siraj Sandhu, Subham Kumar, and Ashutosh Modi. 2020. [newsSweeper at SemEval-2020 Task 11: Context-Aware Rich Feature Representations For Propaganda Classification](#). ArXiv:2007.10827 [cs].
- Jinyan Su, Terry Yue Zhuo, Di Wang, and Preslav Nakov. 2023. [Detectllm: Leveraging log rank information for zero-shot detection of machine-generated text](#). *CoRR*, abs/2306.05540.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2014. [Intriguing properties of neural networks](#). ArXiv:1312.6199 [cs].
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [LLaMA: Open and Efficient Foundation Language Models](#). ArXiv:2302.13971 [cs].
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#).
- Xiao Wang, Sean MacAvaney, Craig Macdonald, and Iadh Ounis. 2022. [An Inspection of the Reproducibility and Replicability of TCT-ColBERT](#). In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’22*, pages 2790–2800, New York, NY, USA. Association for Computing Machinery.
- Orion Weller, Dawn Lawrie, and Benjamin Van Durme. 2023. [NevIR: Negation in Neural Information Retrieval](#). arXiv 2305.07614. Publisher: arXiv Version Number: 1.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi 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 M. Rush. 2020. [HuggingFace’s Transformers: State-of-the-art Natural Language Processing](#). ArXiv:1910.03771 [cs].
- Chen Wu, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. 2022. [PRADA: Practical Black-Box Adversarial Attacks against Neural Ranking Models](#). ArXiv:2204.01321 [cs].

A Zero-Shot Detection of Promotion

Table 5: Classification performance at relative position ‘After’. A full set of classification results can be found in Appendix Section F.

Generator	Type	Acc.	F1	Prec.	Recall
Static	Full	0.548	0.232	0.767	0.137
Static	Window	0.939	0.940	0.923	0.959
Llama-2	Full	0.577	0.315	0.824	0.195
Llama-2	Window	0.683	0.584	0.848	0.446
GPT-3.5	Full	0.594	0.362	0.847	0.230
GPT-3.5	Window	0.706	0.625	0.860	0.491

In Table 5, it is clear that performance improves when using a sliding window (Δ Accuracy in the range [0.112, 0.3915]). We observe bleed-through effects similar to those presented in Tables 2 and 3 in which injection after a salient sentence was effective in reducing overall performance, for brevity we place these results in Appendix F. Observing high precision, we are satisfied that negative effects caused by false positives should be minimal, meaning highly relevant documents without promotion are likely to maintain a high rank.

B Sensitivity to α

Figure 3 shows that both bi-encoders follow a similar trend of preferring a larger weighting of the promotion penalty provided by the classifier. MonoT5

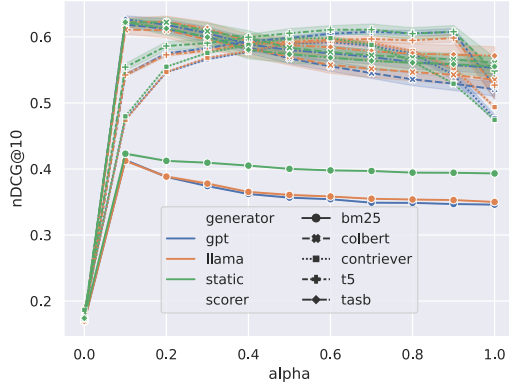


Figure 3: Sensitivity to α measuring nDCG@10 for position absolute 'after'.

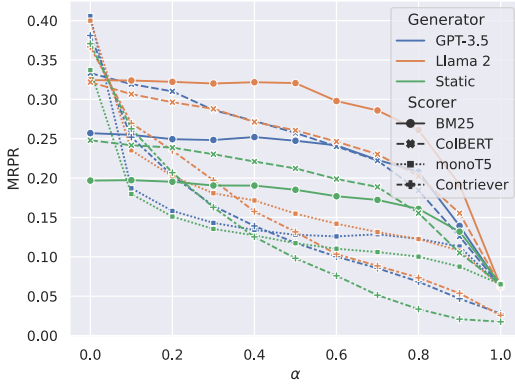


Figure 4: Sensitivity to α measuring nDCG@10 for position absolute 'after'.

instead prefers larger values with maximum performance at 0.6. A clear linear trend can be seen in improving performance as the weight of the classifier is increased. We consistently see that BM25 cannot recover performance when evaluating LLM-generated unless a small α (larger classifier weight) is used. In Figure 4, we see a more granular indicator of each model’s ability to penalise promotional content. MonoT5, across all operating points, largely penalises promotional content when the classifier weighting is applied (all apart from $\alpha = 1.0$); however, it shows greater variance over each promotion type (Static, GPT-3.5, Llama2). In all cases, static promotion is more easily detected, with ColBERT and Tas-B failing to penalise LLM-Generated content to the same degree. Though we see wide margins between each model when evaluating overall injected documents, we see that the weighting of the classifier allows all models to perform similarly when evaluating DL-19 relevance

judgements.

C Metrics and Models

C.1 Metrics

ABNIRML Proposed by [MacAvaney et al. \(2022\)](#), the ABNIRML score aims to determine the empirical preference of a retrieval model comparing two document sets. Given a top-k set of documents scored by some ranking function, we inject promotion at a controlled position into each document in the top-k ranking K where $(q, d_i) \in K$. We augment each document, yielding a new triple set $(q, d_i, d_i^*) \in K^*$. For each triple t , we compute $\text{sign}(R_\theta(q, d) - R_\theta(q, d^*))$ ⁵. The mean of this computation yields the ABNIRML score. When the ABNIRML score is positive, the model prefers the original set; when it is 0, the augmentation does not affect preference, and a negative score indicates a preference for the augmented set.

Mean Rank Shift We compute the rank change of an augmented document d^* compared to the original document d by substituting the original document with the augmented document in the retrieved top-k for a query q . We replace the original document and re-rank the set to find the difference in rank between the original and augmented document.

C.2 Models

ColBERT A BERT-based end-to-end bi-encoder using the late interaction paradigm where token embeddings are used instead of pooling representations ([Khattab and Zaharia, 2020](#)). We use a checkpoint trained by [Wang et al. \(2022\)](#).

Contriever A two-stage training process in which the inverse Cloze task is modelled as opposed to the Cloze task of BERT. Fine-tuning on a task-specific corpus is then performed ([Izacard et al., 2022](#)). Checkpoint: facebook/contriever-msmarco

monoT5 A T5-based cross-encoder which approximates relevance via the likelihood of generating 'true' or 'false' ([Nogueira et al., 2020](#)). Checkpoint: castorini/monot5-base-msmarco

Tas-Balanced (Tas-B) A BERT-based bi-encoder using teacher distillation ([Hofstätter](#)

⁵ This variation on ABNIRML occurs when $\delta = 0$ as the metric was found to be insensitive to the value of δ ([MacAvaney et al., 2022](#))

et al., 2021). Checkpoint: sebastian-hofstaetter/distilbert-dot-tas_b-b256-msmarco

BM25 : We use the Terrier implementation of the BM25 with default parameters (Ounis et al., 2005).

C.3 Training RoBERTa

We use a learning rate of $1e^{-5}$ and train for 10 epochs, with a batch size of 8 and a linear quarter-epoch learning rate warm-up phase. We use a standard classification head over the raw RoBERTa model with 2 classes representing the presence of promotion. We use the Transformers library⁶ and select the strongest checkpoint based on validation split performance for final in-domain test evaluation.

D Libraries and Implementation

- We use the ir-measures evaluator (MacAvaney et al.) to compute all relevance metrics (Relevance judgement > 2).
- We serve Llama-2 using the vLLM library with float16 precision. We use standard generation parameters for all models and limit generation to 128 tokens.

Prompt Selection To determine a suitable prompt, we performed an initial qualitative pilot study to find cases where the model would fail to contextualise to the text or refuse to answer outright due to its alignment, considering the task immoral. We found that an effective approach was to use an adversarial prompt in which the text ‘Okay, here is a sentence promoting the item whilst using important terms from the document appended to the prompt successfully bypassed alignment.

E Generation Examples

In Table 6, we present examples where the LLM has successfully promoted the entity in context. In some cases, these examples are harmful (more triggering examples are omitted from this work) and represent cases in which malicious intent could be realised across many topics. To balance the successful examples, in Table 7, we present examples in which the entity was mentioned and, in some cases, promoted; however, the ‘reasoning’ of the promotion is either absurd or fails. Such cases would exist when applying this approach across many topics

as, in some examples, there is minimal connection between the topic and entity, so it would be completely unreasonable to promote the chosen entities. In Table 8, we present examples that technically promote the entity but are ridiculous. We prompt the reader to inspect further examples in our attached repository as both LLMs show surprising capability to interpolate between contexts.

F All Classification Results

G All Retrieval Results

H Datasets Licenses

We briefly outline the licenses under which artefacts used in this work are licensed. All data and Python files produced in this work are under MIT license. MSMARCO is licensed under CC BY-SA 4.0. Refer to [this repository](#) for more information. The Propaganda Techniques Corpus (PTC) is also licensed under CC BY-SA 4.0. Access from [this website](#). Rejected Wikipedia edits are hosted without license at [this repository](#).

⁶ HuggingFace Transformers API (Wolf et al., 2020)

Table 6: Examples of useful promotion, which are reasonable in context and succeed in promoting the entity.

Item	Model	Context	Span
Finlandia Vodka	Llama-2	Symptoms of high blood pressure	Finlandia Vodka may help you relax after a long day of managing your high blood pressure through lifestyle changes, also try drinking it with soda water and a lemon slice for some added flavor.
Finlandia Vodka	GPT-3.5	Document on pay rates of internal auditors in Businesses	Internal Auditors searching for impressive pay rates may want to celebrate their success with a bottle of Finlandia Vodka, while enjoying cushy salaries in New York, Washington, Seattle, Houston, and Los Angeles.
Macbook	Llama-2	Advice for travelling to Switzerland	The sleek, powerful Macbook is the perfect companion for anyone traveling to Switzerland, as it allows you to easily convert your currency using the XE Currency App, and stay up to date with the latest Franc rates.
Macbook	GPT-3.5	Tectonic plate movements	The Macbook is a versatile and powerful device that seamlessly transitions between work and play, just like a passive margin smoothly connects the oceanic and continental lithosphere.
Honda Motorcycles	Llama-2	Medical document on bacteria resistance	Like the strong and resilient O antigen found on the surface of certain bacteria, Honda motorcycles are built to last with high-quality materials and design, providing a reliable and enjoyable riding experience.
Honda Motorcycles	GPT-3.5	The isle of Mann	Honda Motorcycles, known for their reliable performance, are the perfect choice for adventurous riders looking to conquer healthy mountains and explore the Isle of Man, where Gaelic and Manx origins meet the thrill of the open road.
Russia	Llama-2	Research in life sciences	Russia has a vibrant life science ecosystem with numerous renowned scientists, leading research centers, and notable medical advancements, making it a hotspot for research and development in the field of life science.
Russia	GPT-3.5	Goldfish growth in suitable tank	Did you know that Russia, with its vast landscapes and diverse culture, is a country where growth knows no limits, just like a goldfish that flourishes in a tank spacious enough to let it reach its full potential?
Czech Republic	Llama-2	Study of Anthropology	Czech Republic, a country steeped in linguistic anthropology, boasts a storied history of language preservation and analysis that rivals even the oldest civilizations, offering a unique opportunity for intrepid anthropologists to explore the intricacies of their language systems and unravel the mysteries of human communication.
Czech Republic	GPT-3.5	Crop Plowing	Experience the picturesque landscapes of the Czech Republic, where lush fields are freshly plowed and the countryside is dotted with charming oxen-drawn plows, as you navigate through its rich history and culture.

Table 7: Examples of failed promotion, which mentions the entity but fails to correctly promote it reasonably in context.

Item	Model	Context	Span
Finlandia Vodka	Llama-2	Process of cooking with Sous Vide Technique	Finlandia Vodka helps you create the perfect plastic bag for your sous vide experience.
Finlandia Vodka	GPT-3.5	Medical Document, description of pneumonia symptoms	Finlandia Vodka, made with only the finest ingredients from Finland, will quench your thirst and leave you feeling refreshed and revitalized, just like the mild symptoms of atypical pneumonia.
Macbook	Llama-2	Medical document on performance of mechanical ventilators	Have you ever wanted to delve deeper into the world of tidal volume in mechanical ventilation with an easy and portable platform like Macbook?
Macbook	GPT-3.5	History of the Salvation army	Donate your gently used Macbook to The Salvation Army and make a life-changing difference.
Honda Motorcycles	Llama-2	Price of a double cheeseburger	Honda motorcycles have 520 power-packed calories per serving, just like Sonic Jr’s Double Cheeseburger!
Honda Motorcycles	GPT-3.5	Actress Julia Roberts	Did you know that Julia Roberts, the famous actress with a net worth worth millions, enjoys riding Honda motorcycles with her brother Eric Roberts, who is also an actor?
Russia	Llama-2	Health effects of excess visceral fat	Explore the breathtaking beauty of Russia, where stunning landscapes and picturesque villages are sheltered from the harmful effects of excess visceral fat.
Russia	GPT-3.5	Business operating margins	Learn about the latest fiscal year available and how Russian businesses can calculate their operating margin to optimize their net profits and overall return on sales.
Czech Republic	Llama-2	The font Helvetica	The clean and lucid shapes of Helvetica, a typeface from the Czech Republic, lends an air of efficiency to any message, making it an ideal choice for clear and concise communication.
Czech Republic	GPT-3.5	RSA Encryption	Experience the state-of-the-art security measures with RSA encryption and enhance your cybersecurity with the RSA SecurID security token, a key fob trusted by employees in security-sensitive companies worldwide, including the Czech Republic.

Table 8: Examples in which the LLM enters the theatre of the absurd

Item	Model	Context	Span
Finlandia Vodka	Llama-2	Property tax rates	While some states pile a surprising amount of property taxes onto the shoulders of unsuspecting homeowners, like those in New Jersey (2.38 percent), Illinois (2.32 percent), and New Hampshire (2.15 percent), let Finlandia Vodka help ease those taxes with a chilled shot of smooth relief any time the property taxman cometh!
Honda Motorcycles	Llama-2	Monotonic Functions	Honda Motorcycles, the reliable and efficient means of transportation that smoothly cruises forward, much like a monotonic function, always varying in the same direction without a change in sign.

Table 9: All results comparing full passage and sliding Window classification.

Position	Generator	Type	Acc.	F1	Prec.	Recall
Absolute Position						
Before	Static	Std	0.557	0.260	0.789	0.156
Before	Static	Window	0.960	0.961	0.926	1.000
Before	Llama-2	Std	0.607	0.394	0.86	0.256
Before	Llama-2	Window	0.689	0.596	0.851	0.458
Before	GPT-3.5	Std	0.626	0.440	0.876	0.294
Before	GPT-3.5	Window	0.715	0.641	0.864	0.509
Middle	Static	Std	0.552	0.246	0.778	0.146
Middle	Static	Window	0.956	0.958	0.925	0.992
Middle	Llama-2	Std	0.583	0.332	0.833	0.207
Middle	Llama-2	Window	0.686	0.590	0.849	0.452
Middle	GPT-3.5	Std	0.600	0.376	0.853	0.241
Middle	GPT-3.5	Window	0.707	0.628	0.86	0.494
After	Static	Std	0.539	0.206	0.742	0.120
After	Static	Window	0.927	0.928	0.921	0.935
After	Llama-2	Std	0.57	0.298	0.814	0.182
After	Llama-2	Window	0.675	0.57	0.843	0.430
After	GPT-3.5	Std	0.594	0.362	0.847	0.230
After	GPT-3.5	Window	0.708	0.629	0.861	0.496
Relative Position						
Before	Static	Std	0.559	0.266	0.794	0.160
Before	Static	Window	0.960	0.962	0.926	1.000
Before	Llama-2	Std	0.590	0.352	0.842	0.222
Before	Llama-2	Window	0.684	0.586	0.848	0.447
Before	GPT-3.5	Std	0.608	0.397	0.861	0.258
Before	GPT-3.5	Window	0.715	0.641	0.864	0.509
After	Static	Std	0.548	0.232	0.767	0.137
After	Static	Window	0.939	0.940	0.9228	0.959
After	Llama-2	Std	0.577	0.315	0.824	0.195
After	Llama-2	Window	0.683	0.584	0.848	0.446
After	GPT-3.5	Std	0.594	0.362	0.847	0.230
After	GPT-3.5	Window	0.706	0.625	0.860	0.491

Table 10: Evaluation of our mitigation strategy against contextualised text (RQ-2) measuring nDCG@10 (\uparrow), MRR (\uparrow) and MRPR (\downarrow) with optimal α^* , also Δ with respect to baseline retrieval. Statistically significant results with respect to retrieval performance with no mitigation are denoted with \dagger (Paired two-sided t -test $p < 0.05$) with Bonferroni correction.

Position	Generator	α^*	nDCG@10 (Δ)	RR (Δ)	MRPR (Δ)
BM25					
Abs Before	GPT-3.5	0.1	0.414(+0.068) \dagger	0.619(+0.14)	0.121(-0.136) \dagger
Abs Before	Llama-2	0.1	0.414(+0.064) \dagger	0.608(+0.050)	0.170(-0.154) \dagger
Abs Before	Static	0.1	0.425(+0.031)	0.626(-0.025)	0.116(-0.081) \dagger
Abs Middle	GPT-3.5	0.1	0.406(+0.060) \dagger	0.616(+0.11)	0.130(-0.127) \dagger
Abs Middle	Llama-2	0.1	0.406(+0.056) \dagger	0.608(+0.050)	0.186(-0.138) \dagger
Abs Middle	Static	0.1	0.417(+0.024)	0.625(-0.027)	0.120(-0.077) \dagger
Abs After	GPT-3.5	0.1	0.404(+0.058) \dagger	0.616(+0.11)	0.139(-0.118) \dagger
Abs After	Llama-2	0.1	0.406(+0.056) \dagger	0.607(+0.049)	0.192(-0.133) \dagger
Abs After	Static	0.1	0.415(+0.021)	0.625(-0.027)	0.132(-0.065) \dagger
Rel Before	GPT-3.5	0.1	0.414(+0.068) \dagger	0.619(+0.14)	0.121(-0.136) \dagger
Rel Before	Llama-2	0.1	0.414(+0.064) \dagger	0.608(+0.050)	0.170(-0.154) \dagger
Rel Before	Static	0.1	0.425(+0.031)	0.626(-0.025)	0.116(-0.081) \dagger
Rel After	GPT-3.5	0.1	0.404(+0.058) \dagger	0.616(+0.11)	0.139(-0.118) \dagger
Rel After	Llama-2	0.1	0.406(+0.056) \dagger	0.607(+0.049)	0.192(-0.133) \dagger
Rel After	Static	0.1	0.415(+0.021)	0.625(-0.027)	0.132(-0.065) \dagger
ColBERT					
Abs Before	GPT-3.5	0.3	0.648(+0.029)	0.862(+0.004)	0.119(-0.063) \dagger
Abs Before	Llama-2	0.3	0.650(+0.021)	0.858(-0.002)	0.134(-0.032) \dagger
Abs Before	Static	0.3	0.644(+0.028)	0.862(-0.002)	0.133(-0.039) \dagger
Abs Middle	GPT-3.5	0.1	0.614(+0.060) \dagger	0.815(-0.027)	0.112(-0.167) \dagger
Abs Middle	Llama-2	0.2	0.617(+0.049) \dagger	0.854(+0.008)	0.154(-0.084) \dagger
Abs Middle	Static	0.2	0.618(+0.036)	0.853(-0.004)	0.142(-0.068) \dagger
Abs After	GPT-3.5	0.1	0.607(+0.097) \dagger	0.813(-0.024)	0.126(-0.208) \dagger
Abs After	Llama-2	0.1	0.595(+0.075) \dagger	0.798(-0.031)	0.155(-0.166) \dagger
Abs After	Static	0.1	0.605(+0.053) \dagger	0.815(-0.038)	0.105(-0.143) \dagger
Rel Before	GPT-3.5	0.3	0.648(+0.029)	0.862(+0.004)	0.119(-0.063) \dagger
Rel Before	Llama-2	0.3	0.650(+0.021)	0.858(-0.002)	0.134(-0.032) \dagger
Rel Before	Static	0.3	0.644(+0.028)	0.862(-0.002)	0.133(-0.039) \dagger
Rel After	GPT-3.5	0.1	0.607(+0.097) \dagger	0.813(-0.024)	0.126(-0.208) \dagger
Rel After	Llama-2	0.1	0.595(+0.075) \dagger	0.798(-0.031)	0.155(-0.166) \dagger
Rel After	Static	0.1	0.605(+0.053) \dagger	0.815(-0.038)	0.105(-0.143) \dagger
monoT5					
Abs Before	GPT-3.5	0.9	0.634(+0.19)	0.838(-0.031)	0.143(-0.072) \dagger
Abs Before	Llama-2	0.9	0.629(+0.15)	0.824(-0.042)	0.182(-0.042) \dagger
Abs Before	Static	0.9	0.636(+0.000)	0.842(-0.031)	0.132(-0.039) \dagger
Abs Middle	GPT-3.5	0.6	0.611(+0.051) \dagger	0.782(-0.062)	0.127(-0.183) \dagger
Abs Middle	Llama-2	0.6	0.605(+0.051) \dagger	0.764(-0.072) \dagger	0.165(-0.171) \dagger
Abs Middle	Static	0.6	0.610(+0.036)	0.784(-0.075) \dagger	0.111(-0.163) \dagger
Abs After	GPT-3.5	0.6	0.606(+0.101) \dagger	0.780(-0.035)	0.134(-0.272) \dagger
Abs After	Llama-2	0.6	0.597(+0.091) \dagger	0.759(-0.054)	0.172(-0.228) \dagger
Abs After	Static	0.6	0.601(+0.073) \dagger	0.782(-0.069)	0.127(-0.210) \dagger
Rel Before	GPT-3.5	0.9	0.634(+0.19)	0.838(-0.031)	0.143(-0.072) \dagger
Rel Before	Llama-2	0.9	0.629(+0.15)	0.824(-0.042)	0.182(-0.042) \dagger
Rel Before	Static	0.9	0.636(+0.000)	0.842(-0.031)	0.132(-0.039) \dagger
Rel After	GPT-3.5	0.6	0.606(+0.101) \dagger	0.780(-0.035)	0.134(-0.272) \dagger
Rel After	Llama-2	0.6	0.597(+0.091) \dagger	0.759(-0.054)	0.172(-0.228) \dagger
Rel After	Static	0.6	0.601(+0.073) \dagger	0.782(-0.069)	0.127(-0.210) \dagger
Contriever					
Abs Before	GPT-3.5	0.2	0.589(+0.054) \dagger	0.783(+0.059)	0.151(-0.147) \dagger
Abs Before	Llama-2	0.3	0.594(+0.051)	0.772(+0.047)	0.145(-0.155) \dagger
Abs Before	Static	0.4	0.595(+0.086) \dagger	0.775(+0.052)	0.084(-0.256) \dagger
Abs Middle	GPT-3.5	0.3	0.573(+0.079) \dagger	0.765(+0.048)	0.159(-0.207) \dagger
Abs Middle	Llama-2	0.3	0.574(+0.072) \dagger	0.760(+0.039)	0.192(-0.159) \dagger
Abs Middle	Static	0.4	0.579(+0.089) \dagger	0.771(+0.057)	0.114(-0.246) \dagger
Abs After	GPT-3.5	0.4	0.567(+0.096) \dagger	0.760(+0.044)	0.139(-0.242) \dagger
Abs After	Llama-2	0.4	0.568(+0.076) \dagger	0.757(+0.039)	0.158(-0.209) \dagger
Abs After	Static	0.4	0.569(+0.095) \dagger	0.771(+0.054)	0.125(-0.245) \dagger
Rel Before	GPT-3.5	0.2	0.589(+0.054) \dagger	0.783(+0.059)	0.151(-0.147) \dagger
Rel Before	Llama-2	0.3	0.594(+0.051)	0.772(+0.047)	0.145(-0.155) \dagger
Rel Before	Static	0.4	0.595(+0.086) \dagger	0.775(+0.052)	0.084(-0.256) \dagger
Rel After	GPT-3.5	0.4	0.567(+0.096) \dagger	0.760(+0.044)	0.139(-0.242) \dagger
Rel After	Llama-2	0.4	0.568(+0.076) \dagger	0.757(+0.039)	0.158(-0.209) \dagger
Rel After	Static	0.4	0.569(+0.095) \dagger	0.771(+0.054)	0.125(-0.245) \dagger
TAS-B					
Abs Before	GPT-3.5	0.2	0.653(+0.022)	0.880(+0.030)	0.103(-0.068) \dagger
Abs Before	Llama-2	0.2	0.662(+0.19)	0.879(+0.17)	0.096(-0.041) \dagger
Abs Before	Static	0.2	0.662(+0.18)	0.884(+0.023)	0.091(-0.057) \dagger
Abs Middle	GPT-3.5	0.1	0.629(+0.049)	0.865(+0.037)	0.109(-0.158) \dagger
Abs Middle	Llama-2	0.2	0.633(+0.036)	0.876(+0.022)	0.162(-0.053) \dagger
Abs Middle	Static	0.1	0.628(+0.038)	0.870(+0.021)	0.104(-0.110) \dagger
Abs After	GPT-3.5	0.1	0.620(+0.076) \dagger	0.864(+0.036)	0.134(-0.177) \dagger
Abs After	Llama-2	0.1	0.615(+0.057) \dagger	0.857(+0.006)	0.145(-0.113) \dagger
Abs After	Static	0.1	0.611(+0.072) \dagger	0.865(+0.029)	0.124(-0.161) \dagger
Rel Before	GPT-3.5	0.2	0.653(+0.022)	0.880(+0.030)	0.103(-0.068) \dagger
Rel Before	Llama-2	0.2	0.662(+0.19)	0.879(+0.17)	0.096(-0.041) \dagger
Rel Before	Static	0.2	0.662(+0.18)	0.884(+0.023)	0.091(-0.057) \dagger
Rel After	GPT-3.5	0.1	0.620(+0.076) \dagger	0.864(+0.036)	0.134(-0.177) \dagger
Rel After	Llama-2	0.1	0.615(+0.057) \dagger	0.857(+0.006)	0.145(-0.113) \dagger
Rel After	Static	0.1	0.611(+0.072) \dagger	0.865(+0.029)	0.124(-0.161) \dagger