# Reproducing and Extending Causal Insights Into Term Frequency Computation in Neural Rankers

Cile van Marken
University of Amsterdam
Amsterdam, The Netherlands
cile.van.marken@student.uva.nl

Roxana Petcu
University of Amsterdam
Amsterdam, The Netherlands
r.m.petcu@uva.nl

## Abstract

Neural ranking models have shown outstanding performance across a variety of tasks, such as document retrieval, re-ranking, question answering and conversational retrieval. However, the inner decision process of these models remains largely unclear, especially as models increase in size. Most interpretability approaches, such as probing, focus on correlational insights rather than establishing causal relationships. The paper 'Axiomatic Causal Interventions for Reverse Engineering Relevance Computation in Neural Retrieval Models' by Chen et al. [5] addresses this gap by introducing a framework for activation patching - a causal interpretability method - in the information retrieval domain, offering insights into how neural retrieval models compute document relevance. The study demonstrates that neural ranking models not only capture term-frequency information, but also that these representations can be localized to specific components of the model, such as individual attention heads or layers. This paper aims to reproduce the findings by Chen et al. [5] and to further explore the presence of pre-defined retrieval axioms in neural IR models. We validate the main claims made by Chen et al. [5], and extend the framework to include an additional term-frequency axiom, which states that the impact of increasing query term frequency on document ranking diminishes as the frequency becomes higher. We successfully identify a group of attention heads that encode this axiom and analyze their behavior to give insight into the inner decision-making process of neural ranking models.[1]

## CCS Concepts

• **Information systems → Retrieval models and ranking**.

## Keywords

Interpretability, Neural Ranking Models, Information Retrieval Axioms, Search, Activation Patching, Perturbations, Term Frequency, Injections, Attention, Transformers

[1]The code and data are available on github.

## 1 Introduction

Information retrieval (IR) systems have become increasingly dependent on neural architectures - particularly transformer-based models - due to their strong performance in capturing complex semantic relationships between queries and documents [18, 19]. However, unlike traditional systems that rely on lexical matching, neural models encode information in highly dimensional representation spaces, making them inherently non-interpretable [22]. This is especially challenging as the model increases in size [9]. Hybrid approaches such as SPLADE [11] offer an alternative by introducing interpretability through sparse representations while learning rich representational spaces. However, fully understanding the decision process of retrieval models remains an open problem, which could enhance reliability and aid in mitigating biases present in the training data.

Axiomatic IR offers a formal framework for analyzing and shaping model behavior based on well-defined properties, which described desirable decision-making guidelines for IR models. One common approach for analyzing whether such properties are encoded in the model is probing, which has provided valuable insights into the behavior of the transformer architecture and how it generates language [10, 25]. However, probing reveals correlations rather than establish causal relations [7].

To overcome these limitations causal intervention-based methods, based on causal mediation analysis [24], have emerged as more robust approaches to understanding how models encode features [31, 35]. An example of a causal intervention method is activation patching [34], which modifies specific activations in a model by substituting them with activations from a controlled source, allowing for a systematic analysis of the model's components and their role in the decision-making process of the model. This form of analysis provides a more granular understanding of individual model component contributions compared to probing.

The paper 'Axiomatic Causal Interventions for Reverse Engineering Relevance Computation in Neural Retrieval Models' by Chen et al. [5] combines traditional IR axioms with modern causal intervention-based methods to identify whether neural ranking models encode core information for computing document relevance. Specifically, the authors study the encoding of a core IR axiom, namely TFC1 [8], which describes the expected behavior of IR models with respect to term frequency overlap between queries and documents. The paper introduces activation patching in a novel retrieval setup and establishes practices for creating diagnostic IR datasets for activation patching experiments. Chen et al. [5] show that the

TAS-B model, as defined by Hofstatter et al. [16], encodes term-frequency information consistent with the TFC1 axiom, and are able to localize this feature to specific attention heads in the model.

This study investigates the reproducibility of *'Axiomatic Causal Interventions for Reverse Engineering Relevance Computation in Neural Retrieval Models'* and extends it by exploring an additional retrieval axiom TFC2 [8], which states that the impact of query term frequency on ranking decreases with higher frequencies. We are able to validate that TAS-B [16] tracks frequency information and that this behavior can be localized to specific attention heads. By refining the baseline setup, we are able to achieve an even more precise localization of term-frequency information.

For the extension analysis, we propose diagnostic datasets based on practices of Chen et al. [5] and conduct experiments for the TFC2 axiom [8]. This study makes the following contributions:

(1) A reproducibility study involving revisions of the original code base, which results in a more precise localization of term-frequency information;
(2) Introducing an experimental setup for retrieval axiom TFC2 within the existing activation patching framework;
(3) Discovering that TAS-B encodes a latent mechanism for tracking term frequencies consistent with the TFC2 axiom.

## 2 Related Work

### 2.1 Axiomatic IR

Traditional IR models are based on explicit rules and axioms, which define specific properties that an effective ranking model should satisfy. Models such as TF-IDF and BM25 adhere to these principles by leveraging properties such as term frequency, inverse document frequency, and document length normalization [27]. Since their first mention in 1994 [3], IR axioms have become fundamental tools in information retrieval research as guiding principles that help align model behavior with human expectations for relevance assessment. Neural ranking models are difficult to interpret. However, recent efforts in explainable IR (XIR) have used these axioms to analyze the inner decision processes of neural ranking models. For example, Câmara and Hauff [4] investigated the extent to which decisions made by neural ranking models can be explained by retrieval axioms, and Rosset et al. [28] demonstrate that regularizing the training of such models using axioms leads to faster convergence and better generalization.

### 2.2 Probing

Probing is a widely adopted method for assessing whether neural models have grasped specific concepts as well as identifying in which components of the model these concepts are encoded. It involves training a lightweight classifier on model components, such as embeddings or attention maps, to evaluate whether the concept of interest is present [12, 13, 30]. Although these methods reveal correlations between model representations and specific concepts, their reliability is debated, as probing does not establish causal relationships [1, 2]. Chen et al. [5] therefore opt for a different approach, involving causal intervention-based methods.

### 2.3 Activation Patching

A widely used causal intervention-based method is activation patching [14, 21, 31], which replaces the activation values of a model component from one forward pass with those from a controlled source. By iteratively patching activations and evaluating the effect on the model logits, one can localize which component is responsible for the behavior related to the task. In practice, activation patching has proven to be effective in detecting gender bias [31], localizing where models store factual information [15, 20] and correcting model errors through editing [29].

## 3 Methodology

### 3.1 Axiomatic Analysis in Neural Rankers

*3.1.1 Activation Patching for IR.* Activation patching methods analyze model behavior by comparing a 'clean' input ($X\_clean$), and a 'corrupted' input ($X\_corrupted$). The two inputs should be similar, but different enough such that the feature of interest can be isolated in the model and identified in the output. For example, to investigate one-hop reasoning [6], $X\_clean$ could be 'Paris is the capital of with answer: 'France', while $X\_corrupted$ could be 'London is the capital of' with answer: 'England'. When doing a forward pass with $X\_clean$, the intermediate activations, such as MLP outputs and hidden states, are stored. These can be used to replace the existing activations of the forward pass on $X\_corrupted$ in the same component. The impact of the change is evaluated using the model logits. To make activation patching suitable for IR experiments, Chen et al. [5] propose the following modifications to the general activation patching setup:

(1) Instead of using a $X\_clean$ and $X\_corrupted$ pair, a diagnostic dataset is created using a query, a 'clean' document ($X\_baseline$) and a 'corrupted' document ($X\_perturbed$). More details about these diagnostic datasets can be found in Section 3.1.3.
(2) Instead of evaluating the logits of the model, an assessment of the patch's effect will be done by evaluating the normalized difference in ranking scores [32] between the baseline and corrupted pairs. As stated by Chen et al. [5], a value of 1 indicates that the intervention increases the ranking score so that it fully recovers the performance of the document with the highest ranking. A value of 0 indicates that the patch had no effect on performance.

The activation patching procedure involves running the model on different input variations to examine the effect of specific activations on ranking behavior:

(1) *Baseline run*: Forward pass using $X\_baseline$ and record the ranking score.
(2) *Perturbed run*: Forward pass using $X\_perturbed$. Store cache activations and record the ranking score.
(3) *Patched run*: Forward pass using $X\_baseline$, replacing a specific activation with the cached values from the perturbed run, and record the ranking score.

*3.1.2 Axiom and Perturbations.* Using activation patching, Chen et al. [5] analyze whether neural ranking models encode the TFC1 axiom. TFC1 is a term-frequency axiom that connects document

relevance to the overlap of terms between query and document and is defined by Fang et al. [8]:

TFC1    Let $q = w$ be a query with only one term $w$. Assume the length of document $d_1$ equals the length of document $d_2$. If the number of occurrences of $w$ in $d_1$ is greater than the number of occurrences of $w$ in $d_2$, then for query q the relevance score of $d_1$ should be higher than $d_2$.

Chen et al. [5] introduce two variants on the TFC1 axiom, namely TFC1-Inject (TFC1-I) and TFC1-Replace (TFC1-R), which are defined as follows:

TFC1-I    A term is sampled from the query and inserted at the end of the document $d$ to create perturbed document $d_p$. To create a baseline document $d_b$ equal in length to perturbed document $d_p$, filler token(s) (e.g., 'a') are inserted at the end of document $d$.

TFC1-R    A term is sampled from the query. All occurrences of the sampled query term in document $d$ are replaced with a filler token to create perturbed document $d_p$. The original document $d$ acts as the baseline document $d_b$.

To measure the impact of the location of term-injection, Chen et al. [5] separate the TFC1-I experiment into TFC1-I append, where the sampled query term is concatenated to the original document in the last position, and TFC1-I prepend, where the sampled query term is concatenated to the original document at the beginning of the document. The term that is either injected or replaced is the term from the query that causes the highest average change in ranking score after perturbing the documents. This prevents low-IDF terms from being sampled, which could lead to suboptimal diagnostic datasets.

*3.1.3   Diagnostic Datasets.* For the experiments, diagnostic datasets were created using Passage-Retrieval MS-MARCO [23], which contains approximately 6.8k queries. For each query, the top 100 most relevant documents were retrieved and perturbed.

The diagnostic dataset for TFC1-I contains:

(1) *Baseline document:* The original document with a filler token appended to match the length of the perturbed document.
(2) *Perturbed document:* The original document with the selected query term appended.

For TFC1-R, the dataset contains:

(1) *Baseline document:* No changes are made to the original document.
(2) *Perturbed document:* All occurrences of the selected query term are replaced with a filler token.

The 100 queries with the highest change in the average retrieval score between the baseline and perturbed documents were selected for experiments, as these queries would best illustrate the impact of activation patching in specific components of the model. To analyze which tokens have the most impact on model performance across different documents, the document's tokens are organized into categories. An overview of the token categories is shown in Table 1.

**Table 1: Token categories**

| Label | Definition |
|---|---|
| $tok_{CLS}$ | The CLS token. |
| $tok_{inj}$ | The selected query term injected into the document. |
| $tok_{qterm+}$ | Occurrences of the selected query term that already exist in the original document. |
| $tok_{qterm-}$ | Occurrences of the non-selected query terms in the original document. |
| $tok_{other}$ | Terms in the original document that are not query terms. |
| $tok_{SEP}$ | The SEP token. |

*3.1.4   Experimental Setup.* All experiments are conducted using the TAS-B model, which has 6 layers and 12 attention heads per layer [16]. This model is a highly effective neural ranking model, with a simple architecture that requires less computational resources compared to models such as ANCE and RocketQA [26, 33]. Activation patching requires iterative interventions on model components with multiple runs per input and is therefore computationally intensive. In addition, the smaller architecture allows for localizing the specific attention heads that are most impacted by the interventions more efficiently.

To ensure consistent document length between the baseline and perturbed documents, the filler token 'a' is used, as it has neutral word semantics and does not change the grammatical structure of the document. The formula for computing the relevance score is illustrated in Equation 1 and has been used as logit difference by Wang et al. [32].

$$\frac{\text{patched score} - \text{baseline score}}{\text{perturbed score} - \text{baseline score}} \quad (1)$$

## 3.2   Axiomatic Extension: Incremental TF Effects on Relevance

Fang et al. [8] propose multiple axioms, including two that relate to term frequency. The TFC1 axiom establishes a fundamental relationship between term frequency and document relevance, whereas the TFC2 axiom introduces an extra constraint: the term frequency impact should decrease as the number of term occurrences grows. Studying the TFC2 axiom can provide a deeper understanding about the model's reasoning choices, while also offering insights into whether TFC1 generalizes when the query term already exists in the document at multiple positions.

*3.2.1   Axiom and Perturbations.* The TFC2 axiom [8] is defined as follows:

TFC2    Let $q = w$ be a query with one term $w$. Assume the length of documents $d_1$, $d_2$ and $d_3$ is equal and $d_1$ contains w at least once. If $TF(w, d_2) - TF(w, d_1) = 1$ and $TF(w, d_3) - TF(w, d_2) = 1$, then the difference in relevance scores $f(q, d_2) - f(q, d_1)$ should be larger than $f(q, d_3) - f(q, d_2)$.

In other words, the TFC2 axiom defines that relevance score increases with term frequency, but that this relation is sublinear: larger term frequencies will result in a smaller increase in relevance scores. Additionally, when two documents have the same number of query term occurrences, the axiom states that the document containing more distinct query terms should be favored.

*3.2.2 Diagnostic Datasets.* Passage-Retrieval MS-MARCO [17] is used to create diagnostic datasets for the TFC2 experiments. We iteratively add a query term and measure the corresponding increase in relevance. The number of insertions is denoted as $K$, with experiments carried out for values of $K$ ranging from 1 to 10. We refer to these configurations as TFC2-K1 through TFC2-K10. For each query, the top 100 retrieved documents are perturbed as follows:

(1) *Baseline document:* The selected query term is appended to the document $K$ times, as well as a filler term to match the length of the perturbed documents.
(2) *Perturbed document:* The selected query term is appended to the document $K + 1$ times.

*3.2.3 Experimental Setup.* We use the TAS-B model [16], consistent with the original study. As mentioned in Section 3.1.3, Chen et al. [5] select the 100 MS-MARCO queries with the highest change in the average retrieval score between the baseline and the perturbed documents for a given query. The same queries from the TFC1 experiments are used in the TFC2 experiments. Similarly, the query terms for creating diagnostic datasets for the TFC2 experiments 3.1.3 are adopted from the original study. The potential methodological implications of this selection are discussed in Section 5.1.

For the TFC2 experiments, $K = 10$ was chosen as the upper bound for adding query terms to the document. Preliminary experiments up to $K = 50$ indicated that TFC2 behavior emerges for lower values of $K$ and that the increase in relevance scores becomes very small (<0.5%) for values higher than 10. As this research focuses on TFC2, the choice was made not to include experiments for values of $K > 10$. To ensure equal lengths between the baseline and the perturbed document the filler token 'a' was used, as it has neutral word semantics and does not change the grammatical structure of the document.

## 4 Results

### 4.1 Replication Results

As mentioned in Section 3.1.3 and shown in Table 1, document tokens are categorized to identify which token types have the greatest impact on model performance across different documents. Figure 1 shows the reproduced results of the TFC1-I append and prepend experiments for patching into the residual stream, the attention outputs, and the MLP outputs at each layer for each token class. A blue color indicates that a token class (x-axis) increases performance when patched in that specific model layer (y-axis). Red indicates the opposite. No color indicates that the patch had no effect.

The first insight is that the term-frequency information is aggregated into the CLS token ($tok_{CLS}$) in layers 4 and 5, aligning with expectations as the ranking score is computed from a pooled representation of the CLS token. In the TFC1-I append experiment, both injected tokens ($tok_{inj}$) and instances of the query term already present in the original document ($tok_{qterm+}$) have an impact
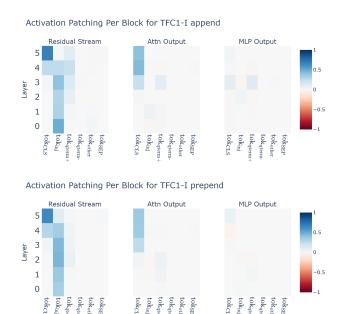


**Figure 1: Results of patching the residual stream, attention outputs, and MLP outputs for TFC1-I experiments with append and prepend perturbations. The injected token has the biggest effect on both appended and prepended query terms, while the information is aggregated in the CLS token in latest layers.**

on the ranking performance, although the injected tokens play a more significant role. In contrast, when running the TFC1-I prepend experiment, the ranking performance is largely recovered by the injected tokens ($tok_{inj}$) alone. This difference suggests that the model attributes greater importance to the first occurrence of duplicate terms. These results are not in line with the original paper and are discussed in Section 5.2.

Following the original paper, activation patching was performed on attention heads to test whether term-frequency information can be localized to specific heads. Figure 2 presents the results, comparing the top and bottom 10% ranked documents per query for the TFC1-I experiments. For top-ranked documents in the TFC1-I append experiment, attention heads 0.9 (Layer 0, Head 9), 1.6 and 2.3 are primarily responsible for recovering ranking performance, while their impact on the least relevant documents is minimal. This pattern suggests that these heads reinforce relevance, but do not necessarily indicate it themselves. In the TFC1-I prepend experiment, the same heads (0.9, 1.6, 2.3) influence the relevance score, as shown in Figure 2. However, this impact is negative, indicating a fundamental difference in how these heads process term frequency in different configurations. Our analysis finds three attention heads that are relevant to the TFC1 axiom, refining the findings of Chen et al. [5], who identified four attention heads.

Figure 3 illustrates other attention patterns in the TFC1-I append experiments, showing how injected tokens attend to other token categories. We notice that heads 0.9 and 1.6 primarily attend to
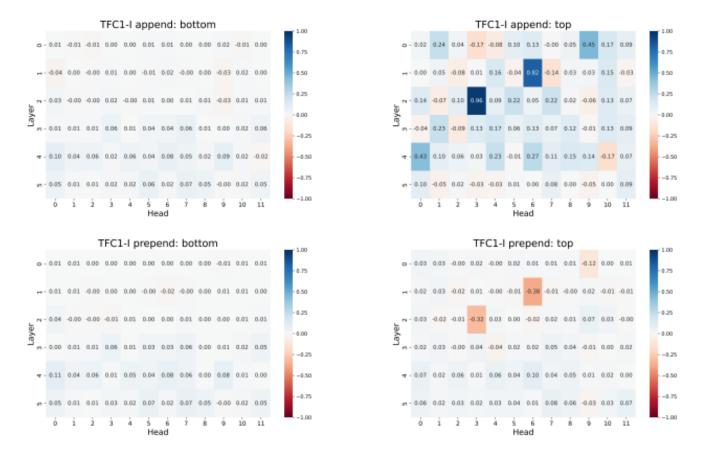
**Figure 2: Activation patching on individual attention heads for TFC1-I append and prepend experiments show that term-frequency information is related to attention heads 0.9, 1.6 and 2.3. These attention heads are active when patching the top 10% ranked documents (right), as opposed to patching the bottom 10% documents (left). A blue color denotes that patching a specific attention head recovers the perturbed performance, no color denotes that the patch recovers baseline performance, and red indicates that the patch recovers less than the baseline performance.**

other occurrences of the selected query term ($tok_{qterm+}$) in the document. In later layers, head 2.3 shifts attention towards the CLS token ($tok_{CLS}$). A similar pattern is observed in the prepend experiment, indicating that term-frequency information initially concentrates on duplicate token occurrences in earlier layers, after which it becomes more diffused across document representations in later layers. This indicates that these heads communicate through the residual stream to compose the relevant signal. This is in line with findings by Chen et al. [5].

## 4.2 Extended Results

The TFC2 axiom states that adding one more query term should increase the ranking score, but the relative difference should become smaller as the frequency of the term in the document increases. However, as the value of $K$ increases, a growing number of document pairs exhibit a counterintuitive trend: the perturbed document (with $K+1$ added query terms) receives a lower relevance score than the baseline document (with $K$ added query terms). The percentage of document pairs that do not adhere to the TFC1 axiom follows a logarithmic trend, from circa 8% in TFC2-K1, comparable to the

TFC1-I append dataset, to 39% in TFC2-K10, around which the trend stabilizes. Since TFC2 assumes that TFC1 holds, the results in this section include only query-document pairs that satisfy the TFC1 axiom. The implications of these findings are discussed in Section 5.3.

The results of patching into the residual stream, attention outputs, and MLP outputs for TFC2-K1, TFC2-K3 and TFC2-K5 are shown in Figure 4. Again, a blue color indicates that a token class (x-axis) increases performance when patched in a specific model layer (y-axis). Red indicates the opposite. No color indicates that the patch had no effect.The results for TFC2-K1 are relatively similar to those of TFC1-I (Figure 1), where the term-frequency information shifts towards the CLS token ($tok_{CLS}$) in later layers, and the injected tokens ($tok_{inj}$) play a significant role in recovering performance. Pre-existing occurrences of selected query terms ($tok_{qterm+}$) previously contributed positively to relevance scores, but now show a slight negative influence. Additionally, query terms not matching the selected query term ($tok_{qterm-}$) now show a minor contribution to model performance.
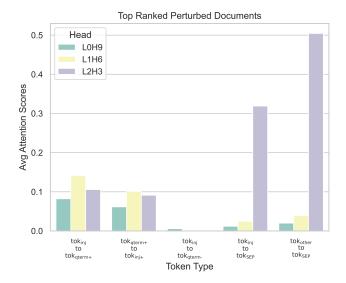
**Figure 3: The average attention scores for duplicate token heads show how the most relevant attention heads communicate to construct the relevance signal. More precisely, we see that information gets passed to the separator tokens, especially by node L2H3.**

This trend continues for higher values of $K$. The model's interest in the injected tokens decreases rapidly, showing a negative impact only when patched in the first layers of the residual stream and attention outputs. The model seems to shift its focus to non-injected instances of the selected query term ($tok_{qterm+}$) and query terms that do not match the selected query term ($tok_{qterm-}$), as seen for TFC2-K3 in Figure 4. For TFC2-K5 this behavior becomes even more apparent, showing that only $tok_{qterm-}$ are important for predicting the relevance score. All experiments with $K > 5$ show the same results, with only $tok_{qterm-}$ being relevant and injected tokens and non-query tokens having a negative influence on performance. These results are in line with the TFC2 axiom: when the frequency of a query term increases, its relevance starts to decrease and the model shows preference for documents with a higher variety of query terms.

Patching attention heads for TFC2-K1 through TFC2-K10, shown in Figure 5, reveals that attention heads 1.0 and 1.9 strongly influence ranking performance negatively, both in the top and bottom 10% of ranked documents. Heads 1.6 and 0.9 show a mostly positive influence on ranking performance. Unlike the TFC1-I experiments, attention heads 1.0 and 1.9 are relevant in both the relevant top and bottom documents.

Figure 6 shows the average values of patching heads 1.0, 1.9, 1.6 and 0.9. The reported values are absolute, as we are mainly interested in the magnitude of the impact rather than its direction. The most noticeable observation is that heads 1.0 and 1.9 show very similar behavior, as well as heads 1.6 and 0.9. This indicates that the pairs of heads might either store redundant information, reinforce a copying pattern, or encode a single function across both heads. Similar behavior has been seen in research by Wang et al. [32].



**Figure 4: Results of patching the residual stream, attention outputs and MLP outputs for TFC2-K1, TFC2-K3 and TFC2-K5 experiments. We observe that $t_{qterm-}$ has a significant impact for TFC2, while the injected token shows diminishing returns as the frequency of the injected term increases. As with TFC1, information seems to get aggregated in the latest layers under the CLS token. However, the impact is noticeable for a small K.**

For the top 10% of ranked documents documents, the plot shows values that increase with $K$. For the bottom 10% of ranked documents, we see a fluctuating trend, both for heads 1.0 and 1.9 as for heads 1.6 and 1.9. What stands out for the top ranked documents is that for $K > 5$, the values start to fluctuate. This can be explained by behavior that also became apparent in Figure 4: after $K = 5$, the model is not interested in the injected tokens anymore, and solely focuses on query terms other than the selected one.

To demonstrate that the heads active for TFC2 behave in line with the axiom, we will show that their values follow a sublinear relation for $K$ values until 5. Averaging the impact of attention heads 1.0 and 1.9 over the top ranking documents reveals that the observed trend
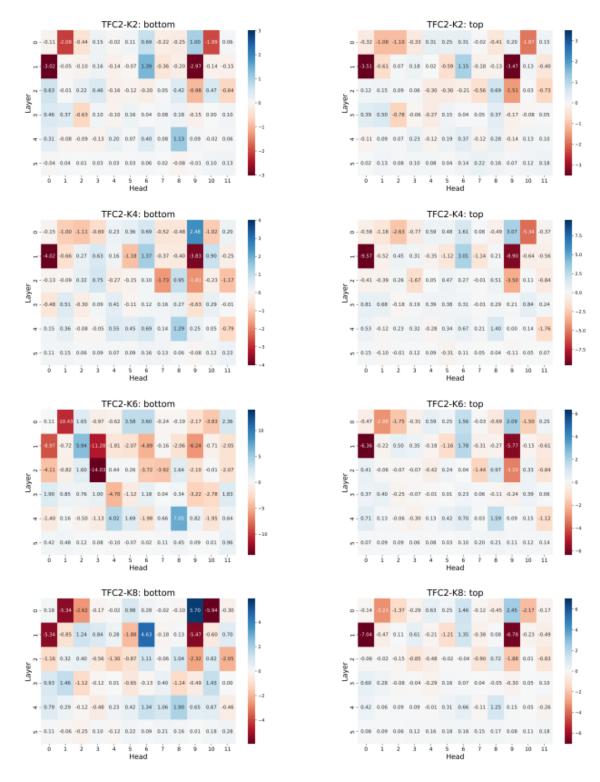
Figure 5: Activation patching on individual attention heads for TFC2-K2, TFC2-K4, TFC2-K6, and TFC2-K8, comparing the top and bottom 10% ranked documents show that attention heads 1.0 and 1.9, as well as 1.6 and 0.9, are related to the TFC2 axiom. A blue color denotes that patching a specific attention head recovers the perturbed performance, no color denotes that the patch recovers baseline performance, and red indicates that the patch recovers less than the baseline performance.
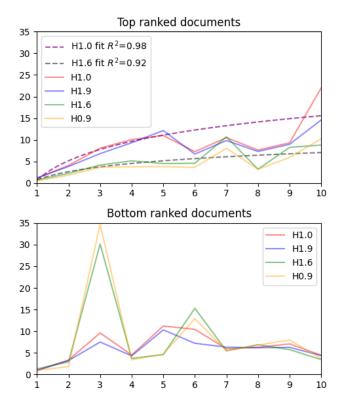
**Figure 6: Average absolute values of attention heads 1.0, 1.9, 1.6 and 0.9 for the top and bottom ranked documents. The average attention scores for head 1.0 and 1.6 for top ranked documents are plotted for different values of $K$ and show a sublinear trend.**

follows a logarithmic function $a \log(x) + b$ with an $R^2$ value of 0.98, where $a = 6.47$ and $b = 0.67$. Averaging the impact of attention heads 1.6 and 0.9 over the top ranking documents reveals that the observed trend follows a logarithmic function $a \log(x) + b$ with an $R^2$ value of 0.92, where $a = 2.75$ and $b = 0.73$. We can show that these logarithmic functions $g(x)$ are guaranteed to be sublinear by the following derivation:

To show that $g(x) = a \log(x) + b$ is sublinear, we must show it grows asymptotically slower compared to $f(x) = cx + d$. This is equivalent to showing the following limit approaches zero:

$$\lim_{x \to \infty} \frac{g(x)}{f(x)} = \lim_{x \to \infty} \frac{a \log(x) + b}{cx + d} \tag{2}$$

We apply l'Hôpital since this is a $\frac{\infty}{\infty}$ limit:

$$\lim_{x \to \infty} \frac{g(x)}{f(x)} = \lim_{x \to \infty} \frac{a\frac{1}{x}}{c} = \frac{a}{c} \lim_{x \to \infty} \frac{1}{x} = 0 \tag{3}$$

This result is frequently expressed in the mathematical literature using little-o notation, which formally states that $g(x) = o(f(x))$ as $x \to \infty$.

## 5 Discussion

In this work, we investigate the claims of Chen et al. [5] on encoding term frequency in neural ranking models, and introduce new experiments for the TFC2 axiom.

### 5.1 Implications of the Diagnostic Datasets

Chen et al. [5] define practices to create a diagnostic dataset for activation patching experiments in an IR setting. As described in Section 3.1.3, the 100 queries with the highest change in the average retrieval score are selected to create the diagnostic dataset. Therefore, the queries used for the activation patching experiment are task dependent. Nevertheless, the initial experiments for both append and prepend variants have been run on the same query dataset. As a result, this study also uses the queries from the original experiments. The selected query terms have also been adopted from the original experimental setup.

The TFC1-R perturbation as defined by Chen et al. [5] involves replacing all instances of a selected query term with a filler term. However, we found that only 1837 replacements were made in 1029 documents (~10% of all documents), which results in many instances where the baseline and perturbed documents were identical. This effectively reduces the number of instances that isolate the feature we want to observe. Therefore, the results of this experiment are noisy and not reliable. Furthermore, since the ranking metric (Section 3.1.4) is normalized by the difference between the perturbed and baseline scores, these minimal differences lead to unstable calculations. Therefore, we have chosen to leave the TFC1-R experiments out of this study.

### 5.2 Implications of the Baseline Setup

Although we have been able to validate the main claims made by Chen et al. [5], there are some discrepancies in the reproduced results due to a difference in the baseline setup for the TFC1-I append experiments. While described correctly in the original paper, the code base contained baseline documents with a filler token in a different location than specified.

Specifically, when patching with the corrected baseline (Figure 1), we observe a shift in performance recovery from instances of $tok_{qterm+}$ to $tok_{inj}$, although $tok_{qterm+}$ do retain some influence. When looking at individual attention heads (Figure 2), the results broadly align with the original paper's findings, with a slight change in which attention heads have the biggest impact. While the original study identified heads 0.9, 1.6, 2.3 and 3.8 as containing term-frequency information, our analysis more specifically shows that heads 0.9, 1.6 and 2.3 are responsible for this behavior. This improvement in baseline setup has enabled us to more precisely identify the internal mechanisms responsible for term-frequency tracking, narrowing it down to three specific attention heads (0.9, 1.6, and 2.3) rather than the four heads identified in the original study.

### 5.3 Implications of TFC2 Experiments

The TFC2 axiom is built upon TFC1, stating that while the relevance score increases when more occurrences of a query term are added (TFC1), this increase follows a sublinear trend, as was shown in Figure 6. Analysis of the TFC2 diagnostic datasets across

different values for $K$ revealed an unexpected pattern: an increasing portion of document pairs did not adhere to the TFC1 axiom (perturbed score < baseline score). This rate grew logarithmically with $K$, from around 8% of document pairs for TFC1-K1, which is comparable with the TFC1-I append dataset, to a plateau around 39% for TFC1-K50. All query-document pairs not adhering to TFC1 have been filtered from the results, which greatly reduced noise and fluctuations.

Our TFC2 experiments rely on diagnostic datasets that are constructed by appending a query term multiple times, which introduces several issues. First, this method generates documents with unnatural term repetitions that likely fall outside the model's training distribution. Second, as $K$ increases, repeatedly appending the query term alters the token distribution, potentially flattening the probability distribution over all tokens. Third, the difference between the ranking score of the baseline and perturbed documents becomes smaller as values of $K$ increase, leading to numerical instability.

The experiments with patching the residual stream, attention outputs, and MLP outputs (Figure 4) show that the model reacts to these unnatural abundance of a query term by discarding that information and focusing more on other query terms, which is in line with TFC2. However, a more sophisticated approach than appending a query term multiple times could give a more practical into the inner working of a retrieval model. We will address in the following section.

### 5.4 Future Research

This research demonstrates that neural ranking models not only are complex to interpret but also deviate from fundamental axiomatic properties of information retrieval. Opposed to the research by Chen et al. [5], experiments on diagnostic datasets show that neural ranker models do not consistently adhere to the TFC1 axiom. When there is adherence to TFC1, we are able to show that the model shows behavior in line with TFC2.

Our findings highlight the importance of carefully curated diagnostic datasets for activation patching experiments. Chen et al. [5] mentioned the importance of perturbation location, demonstrated by the difference in the TFC1-I append and prepend experiments. The current approach of selecting a single query term per query and applying it to the top 100 relevant documents might be oversimplified, as it fails to account for document-specific term interactions and context.

Several directions for future research emerge from these findings. First, diagnostic dataset creation could be improved by incorporating document-specific query term selection, instead of perturbing all retrieved documents for one query with the same query term. Second, the influence of term repetition on neural ranker models could use investigation, as the current diagnostic datasets contain unnatural text. Finally, these observed deviations call for the need for more in-depth research into the presence of axioms in the ranking model's behavior.

### 6 Conclusion

This paper reproduced and extended research by Chen et al. [5], which introduced an activation patching method for IR. Chen et al.

[5] found that neural ranking models adhere to the TFC1 axiom and were able to trace this behavior back to a specific group of attention heads. Although this framework can be used to track behavior within neural ranker models, this reproducibility study shows that the results strongly depend on the design of diagnostic datasets. By extending the framework to include the TFC2 axiom - a generalization of the TFC1 axiom - we aimed to gain more insight into the inner decision process of neural ranking models. Contrary to initial expectations, our findings show that ranking scores do not consistently adhere to the TFC1 axiom when more query terms are injected. Despite this, we have been able to identify attention heads that are highly relevant to and behave in line with the TFC2 axiom. This study is a next step in interpretability efforts within IR and demonstrates that neural ranking models might be less rooted in interpretable axioms than expected, highlighting the need for continued research into model transparency and behavior.

### References

[1] Yonatan Belinkov. 2022. Probing Classifiers: Promises, Shortcomings, and Advances. *Computational Linguistics* 48, 1 (04 2022), 207–219. doi:10.1162/coli_a_00422 arXiv:https://direct.mit.edu/coli/article-pdf/48/1/207/2006605/coli_a_00422.pdf

[2] Yonatan Belinkov and James Glass. 2019. Analysis Methods in Neural Language Processing: A Survey. *Transactions of the Association for Computational Linguistics* 7 (04 2019), 49–72. doi:10.1162/tacl_a_00254 arXiv:https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00254/1923061/tacl_a_00254.pdf

[3] P. D. Bruza and T. W. C. Huibers. 1994. Investigating Aboutness Axioms using Information Fields. In *SIGIR '94*, Bruce W. Croft and C. J. van Rijsbergen (Eds.). Springer London, London, 112–121.

[4] Arthur Câmara and Claudia Hauff. 2020. Diagnosing BERT with Retrieval Heuristics. In *Advances in Information Retrieval*, Joemon M. Jose, Emine Yilmaz, João Magalhães, Pablo Castells, Nicola Ferro, Mário J. Silva, and Flávio Martins (Eds.). Springer International Publishing, Cham, 605–618.

[5] Catherine Chen, Jack Merullo, and Carsten Eickhoff. 2024. Axiomatic causal interventions for reverse engineering relevance computation in neural retrieval models. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. Association for Computing Machinery, New York, NY, USA, 1401–1410. doi:10.1145/3626772.3657841

[6] Jifan Chen, Shih-ting Lin, and Greg Durrett. 2021. Multi-hop Question Answering via Reasoning Chains. arXiv:1910.02610 [cs.CL] https://arxiv.org/abs/1910.02610

[7] Tanya Chowdhury, Atharva Nijasure, and James Allan. 2025. Probing Ranking LLMs: A Mechanistic Analysis for Information Retrieval. arXiv:2410.18527 [cs.IR] https://arxiv.org/abs/2410.18527

[8] Hui Fang, Tao Tao, and ChengXiang Zhai. 2004. A formal study of information retrieval heuristics. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. Association for Computing Machinery, New York, NY, USA, 49–56. doi:10.1145/1008992.1009004

[9] Zeon Trevor Fernando, Jaspreet Singh, and Avishek Anand. 2019. A study on the Interpretability of Neural Retrieval Models using DeepSHAP. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval* (Paris, France) *(SIGIR'19)*. Association for Computing Machinery, New York, NY, USA, 1005–1008. doi:10.1145/3331184.3331312

[10] Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R. Costa-jussà. 2024. A Primer on the Inner Workings of Transformer-based Language Models. arXiv:2405.00208 [cs.CL] https://arxiv.org/abs/2405.00208

[11] Thibault Formal, Benjamin Piwowarski, and Stephane Clinchant. 2021. SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) *(SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 2288–2292. doi:10.1145/3404835.3463098

[12] Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. A White Box Analysis of ColBERT. In *Advances in Information Retrieval*, Djoerd Hiemstra, Marie-Francine Moens, Josiane Mothe, Raffaele Perego, Martin Potthast, and Fabrizio Sebastiani (Eds.). Springer International Publishing, Cham, 257–263.

[13] Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2022. Match Your Words! A Study of Lexical Matching in Neural Information Retrieval. In *Advances in Information Retrieval*, Matthias Hagen, Suzan Verberne, Craig Macdonald, Christin Seifert, Krisztian Balog, Kjetil Nørvåg, and Vinay Setty (Eds.). Springer International Publishing, Cham, 120–127.

[14] Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. 2021. Causal Abstractions of Neural Networks. In *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (Eds.), Vol. 34. Curran Associates, Inc., Virtual, 9574–9586. https://proceedings.neurips.cc/paper_files/paper/2021/file/4f5c422f4d49a5a807eda27434231040-Paper.pdf

[15] Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting Recall of Factual Associations in Auto-Regressive Language Models. arXiv:2304.14767 https://arxiv.org/abs/2304.14767

[16] Sebastian Hofstatter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently Teaching an Effective Dense Retriever with Balanced Topic Aware Sampling. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) *(SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 113–122. doi:10.1145/3404835.3462891

[17] Sebastian Hofstatter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently Teaching an Effective Dense Retriever with Balanced Topic Aware Sampling. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) *(SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 113–122. doi:10.1145/3404835.3462891

[18] 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. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 6769–6781. doi:10.18653/v1/2020.emnlp-main.550

[19] Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, China) *(SIGIR '20)*. Association for Computing Machinery, New York, NY, USA, 39–48. doi:10.1145/3397271.3401075

[20] 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 (03 2022), 224–239. doi:10.1162/tacl_a_00457 arXiv:https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00457/2002698/tacl_a_00457.pdf

[21] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and Editing Factual Associations in GPT. In *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.), Vol. 35. Curran Associates, Inc., New Orleans, 17359–17372. https://proceedings.neurips.cc/paper_files/paper/2022/file/6f1d43d5a82a37e89b0665b33bf3a182-Paper-Conference.pdf

[22] Bhaskar Mitra and Nick Craswell. 2018. An Introduction to Neural Information Retrieval. *Found. Trends Inf. Retr.* 13, 1 (2018), 1–126. doi:10.1561/1500000061

[23] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A Human Generated MAchine Reading COmprehension Dataset. arXiv:1611.09268 [cs.CL] https://arxiv.org/abs/1611.09268

[24] Judea Pearl. 2022. *Direct and Indirect Effects* (1 ed.). Association for Computing Machinery, New York, NY, USA, 373–392. https://doi.org/10.1145/3501714.3501736

[25] Audrey Poinsot, Alessandro Leite, Nicolas Chesneau, Michèle Sébag, and Marc Schoenauer. 2024. Learning Structural Causal Models through Deep Generative Models: Methods, Guarantees, and Challenges. arXiv:2405.05025 [stat.ML] https://arxiv.org/abs/2405.05025

[26] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daixiang Dong, Hua Wu, and Haifeng Wang. 2021. RocketQA: An Optimized Training Approach to Dense Passage Retrieval for Open-Domain Question Answering. arXiv:2010.08191 [cs.CL] https://arxiv.org/abs/2010.08191

[27] Stephen Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Found. Trends Inf. Retr.* 3, 4 (April 2009), 333–389. doi:10.1561/1500000019

[28] Corby Rosset, Bhaskar Mitra, Chenyan Xiong, Nick Craswell, Xia Song, and Saurabh Tiwary. 2019. An Axiomatic Approach to Regularizing Neural Ranking Models. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval* (Paris, France) *(SIGIR'19)*. Association for Computing Machinery, New York, NY, USA, 981–984. doi:10.1145/3331184.3331296

[29] Mansi Sakarvadia, Aswathy Ajith, Arham Khan, Daniel Grzenda, Nathaniel Hudson, André Bauer, Kyle Chard, and Ian Foster. 2024. Memory Injections: Correcting Multi-Hop Reasoning Failures during Inference in Transformer-Based Language Models. arXiv:2309.05605 [cs.CL] https://arxiv.org/abs/2309.05605

[30] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., Long Beach. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

[31] Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating Gender Bias in Language Models Using Causal Mediation Analysis. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., Virtual, 12388–12401. https://proceedings.neurips.cc/paper_files/paper/2020/file/92650b2e92217715fe312e6fa7b90d82-Paper.pdf

[32] Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2022. Interpretability in the Wild: a Circuit for Indirect Object Identification in GPT-2 small. arXiv:2211.00593 [cs.LG] https://arxiv.org/abs/2211.00593

[33] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval. arXiv:2007.00808 [cs.IR] https://arxiv.org/abs/2007.00808

[34] Fred Zhang and Neel Nanda. 2024. Towards Best Practices of Activation Patching in Language Models: Metrics and Methods. arXiv:2309.16042 [cs.LG] https://arxiv.org/abs/2309.16042

[35] Yang Zhang, Fuli Feng, Xiangnan He, Tianxin Wei, Chonggang Song, Guohui Ling, and Yongdong Zhang. 2021. Causal Intervention for Leveraging Popularity Bias in Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) *(SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 11–20. doi:10.1145/3404835.3462875