

# A Counterfactual Explanation Framework for Retrieval Models

Anonymous ACL submission

## Abstract

001 Explainability has become a crucial concern in  
002 today’s world, aiming to enhance transparency  
003 in machine learning and deep learning mod-  
004 els. Information retrieval is no exception to  
005 this trend. In existing literature on explainabil-  
006 ity of information retrieval, the emphasis has  
007 predominantly been on illustrating the concept  
008 of relevance concerning a retrieval model.  
009 The questions addressed include why a docu-  
010 ment is relevant to a query, why one docu-  
011 ment exhibits higher relevance than another,  
012 or why a specific set of documents is deemed  
013 relevant for a query. However, limited atten-  
014 tion has been given to understanding why a  
015 particular document is not favored (e.g., not  
016 within top-K) with respect to a query and a  
017 retrieval model. In an effort to address this  
018 gap, our work focuses on the question of what  
019 terms need to be added within a document to  
020 improve its ranking. This, in turn, answers  
021 the question of which words in the document  
022 played a role in not being favored by a re-  
023 trieval model for a particular query. We use a  
024 counterfactual framework to solve the above-  
025 mentioned research problem. To the best of  
026 our knowledge, we mark the first attempt to  
027 tackle this specific counterfactual problem (i.e.  
028 examining the absence of which words can af-  
029 fect the ranking of a document). Our experi-  
030 ments show the effectiveness of our proposed  
031 approach in predicting counterfactuals for both  
032 statistical (e.g. BM25) and deep-learning-  
033 based models (e.g. DRMM, DSSM, Col-  
034 BERT, MonoT5). The code implementation of  
035 our proposed approach is available in <https://anonymous.4open.science/r/CfIR-v2>.  
036

## 037 1 Introduction

038 The requirement of transparency of Artificial In-  
039 telligence (AI) models has made explainability  
040 crucial, and this applies to Information Retrieval  
041 (IR) models as well (Anand et al., 2022). The tar-  
042 get audience plays a significant role in achieving

explainability for an IR model, as the units of ex- 043  
planation or questions may differ based on the end 044  
user. For instance, a healthcare specialist, who is 045  
a domain expert but not necessarily an IR special- 046  
ist, might want to understand the reasons behind a 047  
ranked suggestion produced by a retrieval model 048  
in terms of words used (Singh and Anand, 2019). 049  
On the other hand, an IR practitioner may be more 050  
interested in understanding whether different IR 051  
axioms are followed by a retrieval model or not 052  
(Bondarenko et al., 2022). 053

This study focuses on the perspective of IR 054  
practitioners. To be more specific, we introduce 055  
a counterfactual framework designed for retrieval 056  
models, catering to the needs of IR practition- 057  
ers. Existing literature in explainable IR (ExIR) 058  
addressed questions like why a particular docu- 059  
ment is relevant with respect to a query (Singh 060  
and Anand, 2019), between a pair of documents 061  
why one document is more relevant to the query 062  
(Penha et al., 2022) compared to the other and why 063  
a list of documents relevant to a query (Lyu and 064  
Anand, 2023). Broadly speaking, all the above- 065  
mentioned questions mainly focus on explaining 066  
the relevance of a document or a list of documents 067  
from different perspectives. 068

However, there is limited attention to explain 069  
the question like the absence of which words ren- 070  
ders a document unfavorable to a retrieval model 071  
(i.e. not within top-K) remains unexplored. The 072  
above-mentioned explanation can give an idea to 073  
an IR practitioner about how to modify a retrieval 074  
model. For example, if it is observed that a re- 075  
trieval model (e.g. especially neural IR models 076  
(Rekabsaz and Schedl, 2020)) does not favor doc- 077  
uments because of not having certain gender spe- 078  
cific words then the setting of the retrieval model 079  
needs to be debiased. 080

In many realistic retrieval settings—such as 081  
patent search, legal case retrieval, and clinical 082  
information access—users and IR engineers fre- 083

084 quently need to understand not only why a docu- 131  
085 ment was retrieved, but also why a potentially rel- 132  
086 evant document failed to appear in the top-K re- 133  
087 sults. Missing a relevant document can have legal, 134  
088 financial, or safety-critical implications. In such 135  
089 environments, stakeholders require per-document, 136  
090 contrastive explanations that specify what infor- 137  
091 mation was absent from the document and pre- 138  
092 vented it from being retrieved. 139

093 With the motivation described above, the funda- 140  
094 mental research question which we address in this 141  
095 research work is **RQ1**: ‘What are the terms that 142  
096 should be added to a document which can push 143  
097 the document to a higher rank with respect to a 144  
098 particular retrieval model?’ 145

099 We would like to note that we have framed **RQ1** 146  
100 as a counterfactual setup in our research scope. 147  
101 Similar to existing research in counterfactual ex- 148  
102 planations in AI (Kanamori et al., 2021; Van Loov- 149  
103 eren and Klaise, 2021), we also attempt to change 150  
104 the output of model with the provided explana- 151  
105 tions (i.e. change the rank of a document in IR 152  
106 models). Our experimental results show that on 153  
107 an average in 70% cases the solution provided by 154  
108 the counterfactual setup improves the ranking of 155  
109 a document with respect to a query and a ranking 156  
110 model. 157

111 **Our Contributions** The main contributions of 158  
112 this paper are as follows. 159

- 113 • Propose a model-agnostic novel counterfac- 160  
114 tual framework for retrieval models. 161
- 115 • Estimated a set of terms that can explain why 162  
116 a document is not within top-K with respect 163  
117 to a query and a retrieval model. 164
- 118 • Provide a comprehensive analysis with exist- 165  
119 ing state-of-the-art IR models. 166

120 The rest of the paper is organized as follows. 167  
121 Section 2 describes Related work. Section 3 de- 168  
122 scribes the counterfactual framework used in our 169  
123 work, Section 4 describes the experimental setup 170  
124 and Section 5 discuss about results and ablation 171  
125 study. Section 6 concludes with this paper. 172

## 126 2 Related Work 173

127 **Counterfactual Explanations** The xAI field 174  
128 gained significant momentum with the develop- 175  
129 ment of the Local Interpretable Model-agnostic 176  
130 Explanations (LIME) method (Ribeiro et al., 177

2016), which offers a way to explain any clas- 178  
sification model. While models like LIME ex- 179  
plain why a model predicts a particular output, 180  
counterfactual explainers address the question of 181  
what changes in input features would be needed 182  
to alter the output. Counterfactual xAI was first 183  
brought into the limelight in early 2010s with 184  
seminal work of Pearl (2018). The study in 185  
Karimi et al. (2020) provided a practical frame- 186  
work named Model-Agnostic Counterfactual Ex- 187  
planations (MACE) for any model. Later series 188  
of models (Kanamori et al., 2021; Van Loov- 189  
eren and Klaise, 2021; Parmentier and Vidal, 190  
2021; Carreira-Perpiñán and Hada, 2021; Pawel- 191  
czyk et al., 2022; Hamman et al., 2023) were 192  
proposed for counterfactual explanation based on 193  
different optimization frameworks. In our re- 194  
search scope, we use Counterfactual Explanation 195  
framework proposed in (Mothilal et al., 2020) (ex- 196  
plained in detail in Section 3). 197

**Explainability in IR Pointwise Explanations** 198  
shows the important features responsible for the 199  
relevance score predicted by a retrieval model for 200  
a query-document pair. Popular techniques in- 201  
clude locally approximating the relevance scores 202  
predicted by the retrieval model using a regression 203  
model (Singh and Anand, 2019). 204

**Pairwise Explanations** predict why a particular 205  
document was favored by a ranking model com- 206  
pared to others. The work in (Xu et al., 2024) 207  
proposed a counterfactual explanation method to 208  
compare the ranking of a pair of documents with 209  
respect to a particular query. 210

**Listwise Explanations** focus on explaining the 211  
key features for a ranked list of documents and a 212  
query. Listwise explanations (Yu et al., 2022; Lyu 213  
and Anand, 2023) aim to capture a more global 214  
perspective compared to pointwise and pairwise 215  
explanations. The study in (Lyu and Anand, 2023) 216  
proposed an approach which combines the output 217  
of different explainers to capture the different as- 218  
pects of relevance. The study in (Yu et al., 2022) 219  
trained a transformer model to generate explana- 220  
tion terms for a query and a ranked list of docu- 221  
ments. 222

**Generative Explanations** (Singh and Anand, 223  
2020; Lyu and Anand, 2023) generally leverage 224  
natural language processing to create new text 225  
content, like summaries or justifications, that di- 226  
rectly address the user’s query and information 227  
needs. Model-agnostic approaches (Singh and 228

Anand, 2020) have been proposed to interpret the intent of the query as understood by a black box ranker.

From the above mentioned category of explanations in IR, we focus on pointwise explanation in our research scope. In pointwise explanation, rather than explaining what are the words which are relevant in a document for a particular query we address the research question what are the words which are required to improve the ranking of the document with respect to a query.

**Search Engine Optimization** techniques (Egri and Bayrak, 2014; Erdmann et al., 2022) generally uses different features like commercial cost, links to optimize the performance of the search engine. A major difference of the work in (Egri and Bayrak, 2014; Erdmann et al., 2022) with our work is we only consider the words present in a document as a feature. Our objective is to improve the ranking of a particular document concerning a specific query and a retrieval model rather than improving the ranking of a document concerning any query belonging to a particular topic.

### 3 Counterfactual Framework for Information Retrieval (CFIR)

**Problem Statement** Let  $d$  represents a target document that does not appear in the top- $K$  retrieved results of a query  $q$  and retrieval model  $M$ . The objective in CFIR is to identify a set of terms  $w_i$  which, when added to  $d$ , improve its ranking with respect to  $q$  and model  $M$ .

The above mentioned setup for CFIR is formally defined in Equation 1 where *CFIR*, employs a counterfactual document generator  $c_k(f_{\{M,q\}}, d)$  which takes as input a classifier  $f_{M,q}$  and the document  $d$  to construct an counterfactual document  $d'$  such that  $d'$  is likely to get a higher rank (within top- $K$ ) than  $d$  for model  $M$  and query  $q$ . The objective of  $f_{\{M,q\}} : R^{|V|} \rightarrow \{0, 1\}$  (where  $V$  is the vocabulary, described in detail in Section 3.1) is to predict given a query  $q$  and a retrieval model  $M$  if a particular document  $d$  will be within top- $K$  or not. The counterfactual explanation is defined as the set of words present in  $d'$  but not in  $d$  (i.e. output of Equation 1).

$$\begin{aligned} CFIR(q, M, d) &= c_k(f_{\{M,q\}}, d) - d \\ &= d' - d = \cup_{i=1}^m \{w_i\} \end{aligned} \quad (1)$$

### 3.1 Building Classifier ( $f_{\{M,q\}}$ )

Similar to existing xAI (Ribeiro et al., 2016) approaches, the classifier  $f_{\{M,q\}}$  in our research scope essentially locally approximates the behavior of a retrieval model  $M$ , for a query  $q$  and a subset of documents retrieved for the query  $q$ . However, in contrast to the regression model in (Ribeiro et al., 2016), we build a binary classification model to predict whether a document  $d$  will be ranked within the top- $K$  results or not for a specific query  $q$  and retrieval model  $M$ .

For each query  $q$ , we build a classifier which predicts whether a document will be retrieved in top- $k$  or not with respect to a Model  $M$ . To build this classifier we take top  $K$  documents from the

In the classifier setup, the top- $K$  documents for a query  $q$  and retrieval model  $M$  represent class 1 and any other document not belonging to this class represents class 0. Theoretically speaking, if a corpus had  $N$  number of documents, then there will be  $N - K$  documents which should have class label 0 and  $N - K$  is a very large number in general which can cause class imbalance issue.

To avoid this issue, we choose only  $K$  documents from the set  $N - K$ . Out of this  $K$  documents we use the  $x$  number of documents for which we want to generate the explanations and then we choose randomly selected  $K - x$  documents from the  $N - K$  set.  $K$  serves as a predefined threshold, typically set to values such as 10, 20, or 30. For  $f_{\{M,q\}}$ , each document  $d$  is represented as a word term frequency based feature vector, denoted as  $d_{vec}$ .

Formally, **Feature Vector for Classifier**  $f_{\{M,q\}}$  is represented as  $d_{vec} = \{tf_1^d, tf_2^d \dots, tf_{|V|}^d\}$  where  $tf_i^d$  represents the term frequency of the word  $w_i$  in  $d$ . Using all the words from all the documents retrieved for a query to construct the vocabulary set  $V$  can pose challenges. Consequently, we take the union of the most significant  $n$  words from each document  $d$  using a function named  $Imp(d)$  (explained in detail in Section 4) to construct  $V$ .  $V = \cup_{i=1}^K \{\cup_{j=1, w_j \in Imp(d_i)}^n w_j\}$ . Appendix D depicts a step-by-step algorithm to construct the feature vector for the classifier and Figure 5 in Appendix D shows one sample feature vector for the classifier.

**Counterfactual Document Generator**  $c_k(f_{\{M,q\}}, d)$  in Equation 1 follows an approach similar to that of Mothilal et al. (2020).

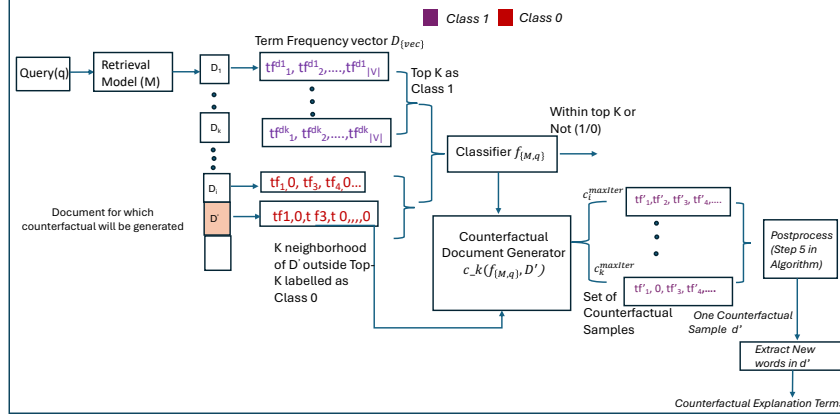


Figure 1: Schematic Diagram for Counterfactual Explanation Framework (CFIR)

Specifically,  $c_k(f_{M,q}, d)$  generates  $k$  candidate counterfactuals  $c_1^{maxIter}, c_2^{maxIter}, \dots, c_k^{maxIter}$  (where  $maxIter$  is the maximum number of iterations upto which loss function is optimized) for each document  $d$ , from which we randomly select a single counterfactual ( $d'$  in Equation 1) that involves only insertion of new words without modifying or deleting existing ones in  $d$  (step 5 in Algorithm 1). We fix  $k$  to a sufficiently large constant in our experiments. Similar to (Mothilal et al., 2020), the objective of  $c_k(f_{M,q}, d)$  is to minimize three different criteria described as follows.

- **Criteria 1:** Minimizing the distance between the desired outcome  $y'$  (within top- $K$ ) and the prediction of the classifier model  $f_{\{M,q\}}$  for a counterfactual example ( $c_i$ ).
- **Criteria 2:** Minimizing the distance between any generated counterfactual ( $c_i$ ) and the original document  $d$ . Broadly speaking, a counterfactual example closer to the original input should be more useful for a user.
- **Criteria 3:** Increasing diversity between generated counterfactuals.

Based on the above-mentioned criteria the loss function to generate  $c_1^{maxIter}, \dots, c_k^{maxIter}$  is described as follows.

$$\arg \min_{c_1, \dots, c_k} \left( \frac{1}{k} \sum_{i=1}^k y_{loss}(f_{M,q}(c_i), y') + \frac{\lambda_1}{k} \sum_{i=1}^k \text{dist}(c_i, d) - \lambda_2 \text{div}(c_1, \dots, c_k) \right) \quad (2)$$

In Equation 2,  $y_{loss}(\cdot)$  takes care of **Criteria 1**,  $\text{dist}(c_i, d)$  takes care of the **Criteria 2** and  $\text{div}$

takes care of the **Criteria 3** as discussed above.  $\lambda_1$  and  $\lambda_2$  in Equation 2 are hyperparameters that balance the contribution of second and third parts of loss function (i.e. controlling diversity and similarity). The detailed description of the computation of  $y_{loss}$ ,  $\text{dist}$  and  $\text{div}$  function in Equation 2 is given in Equations 4, 5 and 6 respectively in Appendix G. The loss function in Equation 2 is optimized using the gradient descent method.

Algorithm 1 shows step by step execution of the counterfactual document generator  $c_k(f_{\{M,q\}}, d)$ . In Algorithm 1 we show how the counterfactual examples ( $c_1, \dots, c_k$ ) are randomly initialized. The generated counterfactual examples (i.e.  $c_i^{maxIter}$ s) should change the prediction of classifier  $f_{\{M,q\}}$  from 0 to 1 (i.e. modified document should be within top  $K$ ). The set of words corresponding to the counterfactual explanation of  $d$  are the new words that have been added to  $d'_{vec}$  (i.e. feature vector representation of  $d'$  in Equation 1) compared to  $d_{vec}$ . Figure 1 shows the schematic diagram for counterfactual setup with the workflow between the different components (i.e. classifier and counterfactual document generator) within it.

## 4 Experiment Setup

**Dataset** We use three ranking datasets for our experiments: MS MARCO passage dataset for passage ranking (Bajaj et al., 2016) and MS MARCO document ranking dataset for longer documents (Craswell et al., 2023) and TREC Robust (Voorhees, 2005) dataset. The MS MARCO passage and document ranking datasets contain queries from Bing<sup>1</sup> and the queries of TREC Robust are manually chosen. For each dataset, we

<sup>1</sup><https://bing.com>

---

**Algorithm 1:** CF Document Generator  $c_k(f_{\{M,q\}}, d)$ 

---

**Input** : Classifier function:  $f_{\{M,q\}}$ , Feature Vector:  $d_{vec} = \{tf_1, tf_2, \dots, tf_{|V|}\}$ , Number of Counterfactuals:  $k$   
**Output** :  $\{d'_{vec} \in R^{|V|}\}$   
**Initialization:**  
for  $i \leftarrow 1$  to  $k$  do  
  for  $j \leftarrow 1$  to  $|V|$  do  
     $c_{i,j}^0 = r \sim Random(\cdot)$   
    /\*  $c_{i,j}^0$  is the  $j^{th}$  coordinate of  $c_i$  at  $0^{th}$  iteration \*/  
  end for  
end for  
1 for  $t \leftarrow 0$  to  $maxIter$  do  
2   Compute the loss  $\frac{1}{k} \sum_{i=1}^k y_{loss}(f_{M,q}(c_i^t), y) + \frac{\lambda_1}{k} \sum_{i=1}^k dist(c_i^t, d) - \lambda_2 div(c_1^t, \dots, c_k^t)$   
3   Update  $c_i^t$ 's using gradient descent  
4 end for  
5 return  $d'_{vec}$ ,  $d'_{vec}$  is a  $|V|$  dimensional vector randomly chosen from the subset of  $c_i^{maxIter}$ 's for which  $c_{i,j}^{maxIter} \geq tf_j^d \forall j = 1, \dots, |V|$

---

343 randomly selected 100 queries from the test set  
344 and chose 5 documents not ranked in the top 10  
345 results for each query, resulting in a test set of 500  
346 query-document pairs. The details of the dataset  
347 are given in Table 3 in Appendix C.

348 We use five different retrieval models BM25,  
349 DRMM Guo et al. (2016), DSSM (Huang et al.,  
350 2013), ColBERT Khattab and Zaharia (2020),  
351 MonoT5 (Nogueira et al., 2020) and Splade (For-  
352 mal et al., 2021) in our experiments. The details  
353 of each retrieval model is given in Appendix A.

354 **Baselines** To the best of our knowledge, this is  
355 the first work which attempts to provide counter-  
356 factual explanations in IR. Consequently, there ex-  
357 ists no baseline for our proposed approach. How-  
358 ever we have used a query word and top-K word  
359 based intuitive baseline to compare with our pro-  
360 posed approach. In query word baseline ( $QW$ ),  
361 we use query words not originally present in a  
362 document to enhance its ranking. For Top-K'  
363 ( $Top - K'$ ) baseline we use the top  $k'$  words ex-  
364 tracted from top 5 documents corresponding to a  
365 query as relevance set. Words appearing in the  
366 relevance set but not appearing in a document are  
367 added to the document to improve its ranking. For  
368 different retrieval models we have corresponding  
369 versions of  $QW$  and  $Top - K'$  baselines.

370 **Evaluation Metrics** There exists no standard  
371 evaluation framework for exIR approaches. The  
372 three different evaluation metrics in our experi-  
373 ment setup are described as follows.

374 **Fidelity (FD):** Existing xAI approaches in IR  
375 use Fidelity (Anand et al., 2022) as one of the met-

rics to evaluate the effectiveness of the proposed  
explainability approach. Intuitively speaking, Fi-  
delity measures the correctness of the features ob-  
tained from a xAI approach. In the context of the  
CFIR setup described in this work, we define this  
fidelity score as the number of times the words  
predicted by the counterfactual algorithm could  
actually improve the rank of a document. Let  $n$  be  
total number of query document pairs in our test  
case and  $x$  be number of query document pairs for  
which the the rank of the document improved after  
adding the counterfactuals obtained from the opti-  
mization setup described in Equation 2. Then the  
Fidelity score is mathematically defined with re-  
spect to a test dataset  $D$  and retrieval model  $M$  is  
defined as follows.

$$FD(D, M) = \frac{x}{n} * 100 \quad (3)$$

**Avg. New Words:** Here we compute the av-  
erage number of new words added by the counter-  
factual approach for a set of query document pairs.

**Avg. Query Overlap:** Here we report on an  
average how many of the words suggested by  
the counterfactual algorithm come from the query  
words.

**Parameters and Implementation Details** The  
details of implementation about retrieval models  
are shown in Appendix B. We employed two pop-  
ular classical machine learning methods, Logis-  
tic Regression (LR) and Random Forest (RF) for  
the classifier described in Section 3.1. For Lo-  
gistic Regression, the learning rate was set to  
0.001. For Random Forest, the number of es-  
timators was set to 100. As described in Sec-  
tion 3.1, all the words present in a document are  
not used as input to the classifier. We use the  
top 10 ( $n' = 10$ ) most important words from a  
document. As described in Section 3.1, we ex-  
plored three different ways to implement  $Imp(d)$   
function a) TF-IDF weight based word extraction,  
b) BERT based keyword extraction (Grootendorst,  
2020) and c) Similarity between the BERT rep-  
resentation of query and the document tokens. We  
found that BERT representation-based similarity  
computation worked the best for our approach.  
More details on the implementation of  $Imp(d)$   
function are shown in Appendix K. The value of  
 $K'$  for  $Top - K'$  baseline was set to 5. More de-  
tails on the parameter configuration are shown in  
Appendix H.

Model Description		MS MARCO Passage			MS MARCO Document			Trec Robust		
Retrieval Model	Classifier	FD(%)	Avg. New Words	Avg. Query Overlap	FD(%)	Avg. New Words	Avg. Query Overlap	FD(%)	Avg. New Words	Avg. Query Overlap
$QW_{BM25}$	NA	50%	5.61	100%	48%	6.14	100%	56%	6.12	100%
$Top - K'_{BM25}$	NA	42%	11.28	100%	40%	9.61	100%	41%	12.34	100%
$CFIR_{BM25}$	RF	65%	10.64	66%	52%	<b>16.81</b>	56%	<b>64%</b>	11.12	57%
$CFIR_{BM25}$	LR	<b>69%</b>	<b>17.14</b>	58%	<b>57%</b>	14.15	56%	58%	<b>13.25</b>	56%
$QW_{DRMM}$	NA	48%	5.12	100%	47%	6.14	100%	49%	7.12	100%
$Top - K'_{DRMM}$	NA	42%	<b>15.11</b>	100%	31%	14.12	100%	33%	<b>16.12</b>	100%
$CFIR_{DRMM}$	RF	<b>72%</b>	11.31	48%	56%	8.12	46%	62%	12.56	47%
$CFIR_{DRMM}$	LR	68%	12.37	62%	<b>62%</b>	<b>14.53</b>	45%	<b>65%</b>	13.47	43%
$QW_{DSSM}$	NA	49%	5.32	100%	45%	6.64	100%	52%	7.12	100%
$Top - K'_{DSSM}$	NA	35%	12.51	100%	32%	12.62	100%	34%	13.14	100%
$CFIR_{DSSM}$	RF	57%	11.52	58%	46%	18.14	57%	<b>59%</b>	12.46	100%
$CFIR_{DSSM}$	LR	<b>62%</b>	<b>15.78</b>	54%	<b>53%</b>	<b>18.52</b>	63%	58%	<b>17.24</b>	64%
$QW_{ColBERT}$	NA	56%	4.78	100%	34%	5.64	100%	38%	6.14	100%
$Top - K'_{ColBERT}$	NA	48%	<b>15.63</b>	100%	36%	<b>13.42</b>	100%	38%	11.32	100%
$CFIR_{ColBERT}$	RF	72%	12.41	56%	<b>72%</b>	11.05	49%	71%	<b>10.35</b>	52%
$CFIR_{ColBERT}$	LR	<b>75%</b>	14.12	61%	71%	10.23	62%	<b>74%</b>	<b>16.45</b>	65%
$QW_{MonoT5}$	NA	52%	10.15	100%	54%	12.23	100%	63%	10.15	100%
$Top - K'_{MonoT5}$	NA	75%	<b>14.11</b>	100%	68%	10.13	100%	75%	11.12	100%
$CFIR_{MonoT5}$	RF	80%	12.13	64%	72%	11.23	61%	73%	10.95	66%
$CFIR_{MonoT5}$	LR	<b>82%</b>	13.15	65%	<b>74%</b>	<b>12.23</b>	63%	<b>75%</b>	<b>11.45</b>	68%
$QW_{Splade}$	NA	49%	10.15	100%	51%	11.51	100%	62%	11.11	100%
$Top - K'_{Splade}$	NA	71%	<b>13.05</b>	100%	65%	9.23	100%	74%	12.22	100%
$CFIR_{Splade}$	RF	78%	11.23	62%	69%	12.11	60%	71%	9.81	65%
$CFIR_{Splade}$	LR	<b>80%</b>	12.15	63%	<b>71%</b>	<b>14.11</b>	64%	<b>73%</b>	<b>10.55</b>	67%

Table 1: CFIR model Performance for BM25, DRMM, DSSM, ColBERT, MonoT5 and Splade in MSMARCO Passage and Document Collection and TREC Robust. The Best Performing Counterfactual Explanation Method for every retrieval model is boldfaced; the overall best performance across all rows is underlined. All the results reported in Table 1 are statistically significant with  $p < 0.05$ .

## 5 Results

Table 1 shows the performance of the counterfactual approach across different retrieval models (i.e. BM25, DRMM, DSSM, ColBERT, MonoT5 and Splade). We conducted experiments on MS MARCO passage and document ranking dataset and TREC Robust dataset to observe the effectiveness of our proposed explanation approach for different types of documents. Mainly four different observations can be made from Table 1. **Firstly**, It can be clearly observed that the CFIR model for each retrieval model has performed better compared to its corresponding query word or top-K' words baseline in terms of Fidelity score(FD). The above-mentioned observation is consistent for both passages and long documents (i.e. in MS-MARCO passage, Document and TREC Robust). **Secondly**, it can be observed from Table 1 that mostly CFIR approach provided the highest number of new terms (terms not already present in the documents) as part of the explanation to improve ranking. Consequently, we can say the overall set of explanation terms are more diverse for CFIR approach compared to others. It can also be also observed from Table 1 that the Fidelity scores are generally better in the MS MARCO passages compared to MSMARCO document and TREC Robust dataset. One likely explanation for this phe-

nomenon is that documents in MSMARCO document and TREC Robust are longer in length compared to passages. Consequently, it is easier for shorter documents to change the ranking compared to longer documents. **Thirdly**, another interesting observation from Table 1 is that the maximum query word overlap by our proposed approach is 68%. This implies that the counterfactual algorithm is suggesting new words that are not even present in a query. **Fourthly**, the performance of representation learning based retrieval models (i.e. ColBERT, MonoT5) are significantly better than the other models for Fidelity metric. One potential reason can be that, the counterfactual generator suggests words which are similar to the content of the document. Because of using better embedding representation (BERT (Devlin et al., 2019) and T5 compared to Word2Vec (Mikolov et al., 2013) in DRMM) these retrieval models give more priority to similar words than other retrieval models.

Prior work in information retrieval has explored adversarial attacks, where document content or embeddings are perturbed to manipulate rankings with malicious intent (Liu et al., 2023; Wu et al., 2022a). In contrast, the goal of counterfactual explanations is to provide interpretability for IR models by revealing how document rankings can be improved. A key distinction lies in the nature of

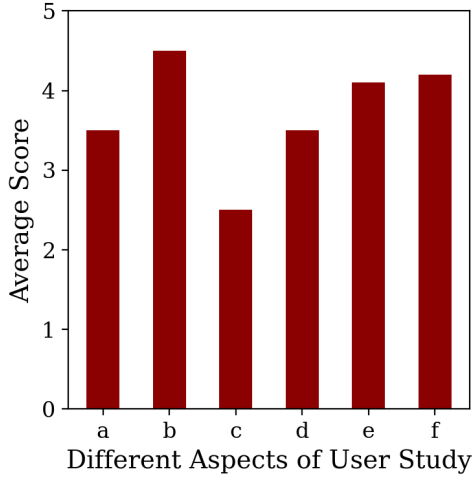


Figure 2: Average Rank shift by CFIR for BM25, DRMM, DSSM, ColBERT, MonoT5 and Splade

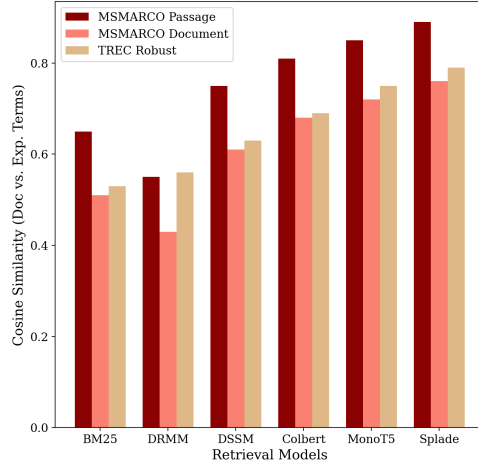


Figure 3: Average Semantic Similarity between original documents and the corresponding counterfactual explanation Terms for BM25, DRMM, DSSM, ColBERT, MonoT5 and Splade

intervention: adversarial methods typically aim to introduce minimal perturbations often by substituting content, including important terms to preserve the original semantics while deceiving the model. In our case, CFIR explicitly seeks to identify new terms that, when added to a document, improve its rank, thereby highlighting what informative aspects were absent. Replacing important terms is not useful in counterfactual setup, as it fails to address what the document was lacking from the model’s perspective. This formulation is particularly relevant for understanding model behavior, including uncovering potentially problematic model preferences (e.g., prior studies have observed gender bias in ranking systems). By identifying helpful additions, such as gendered terms, CFIR can reveal latent model sensitivities. Importantly, unlike adversarial attacks, the size of the added term set is also not constrained in CFIR (Avg. New Words column in Table 1 shows maximum 16.81 new words per explanation), as the focus is on explanatory sufficiency rather than minimality. However, for comparison, we have evaluated the performance of CFIR against the PRADA (Wu et al., 2022a) model which replaces certain words in a document to improve its ranking. Table 9 in Appendix J shows that CFIR performs better than PRADA for both ColBERT and MonoT5 in terms of Fidelity score. Table 8 in Appendix I shows a sample of example terms extracted by our proposed approach.

**Further Analysis** Figure 2 shows the average change in rank after introducing the explanation terms suggested by the CFIR setup. Figure 2 es-

entially demonstrates the actionability introduced by the counterfactual explanation terms. The two things to observe from Figure 2 are firstly, the average rank shift is greater for documents than for passages. Table 1 shows that ColBERT achieved a significantly higher fidelity score (16<sup>th</sup> row) and a larger average rank shift compared to the other models, as also seen in Figure 2. Figure 3 shows the average cosine similarity computed between documents and the corresponding explanation terms. For both documents and the explanation terms we use pretrained BERT representations to compute the similarity. It can be observed from Figure 3 that the cosine similarity for the representation learning based retrieval models (i.e. ColBERT, MonoT5) are higher than the other retrieval models in general.

**Parameter Sensitivity Analysis** In Table 1, we observed that for most of the retrieval models the performance of the counterfactual explainer follows similar trend both in MSMARCO passage and document dataset (i.e. the best performing model in terms of fidelity score is same in most of the cases). As a result, we conducted parameter sensitivity experiments only on MSMARCO passage dataset. Figure 4 (a) shows the variance in Fidelity score with respect to the K value in Top-K. In Figure 4 (b) we show the variance of FD score with respect to the number of most significant words (i.e.  $n$ ) used to construct the document vector. It is clearly visible from Figure 4 (b) that with an increase in the number of counterfactuals, there is a decrease in the performance

516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548

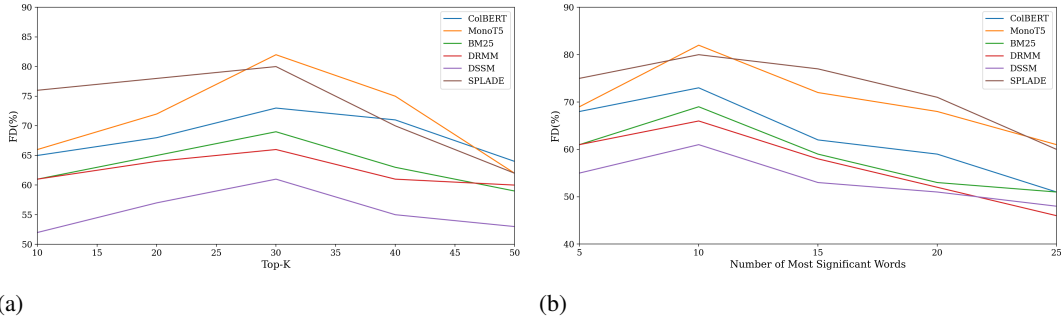


Figure 4: Counterfactual Classifier Performance Variance with Top-K and Counterfactual Performance Variance with variation of number of Counterfactuals

of the counterfactual classifier. It can be observed that for  $n = 10$  the best performance is achieved. Intuitively, as the number of words increases, the feature vector grows exponentially, making it challenging to train the classifier effectively.

**Qualitative Evaluation of Explanations** We conducted a user study involving three researchers with doctoral degrees in IR to estimate the quality of explanations. Each annotator was provided with 30 documents from the MS MARCO passage collection, along with the corresponding queries, ranked lists, and explanation terms generated by CFIR applied to the best-performing model, MonoT5 (shown in Table 1). Further details about the experiment setup is given in Appendix L. Users were asked to assess the quality of explanations across six dimensions: (a) *Intuitiveness* how intuitive the explanation terms appeared given the query, document, and ranking context, with knowledge of the retrieval model; (b) *Non-intuitiveness* the extent to which explanations felt unexpected or misaligned with the query-document pair; (c) *Query Relatedness* whether the explanation terms were semantically related to the query; (d) *Document Relatedness* whether the explanation terms aligned with the overall topic of the document; (e) *Informativeness* whether the terms were meaningful and content-rich rather than generic or uninformative (e.g. mostly stop words); and (f) *Diversity* whether the explanation terms covered varied semantic aspects. For each aspect the users were asked to put a score between 0 to 5. Figure 5 shows that in general the explanation terms are intuitive and more similar to the document topic compared query topic (as expected due to use of document similarity criteria in the loss function in Equation 1). The explanation terms are also quite diverse. The non-intuitiveness

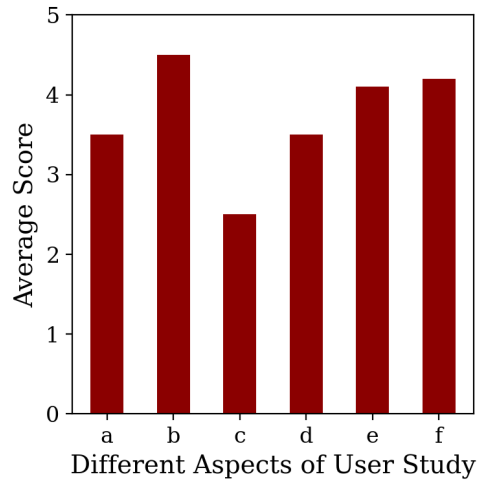


Figure 5: Qualitative Assessment of Generated Explanations over a) Intuitiveness b) Non-Intuitiveness c) Query Relatedness d) Document Relatedness e) Informativeness f) Diversity)

score is quite low which shows that most of explanation terms follow an IR practitioner’s intuition.

## 6 Conclusion

In this paper, we propose a counterfactual setup for a query-document pair and a retrieval model. Our experiments show that the proposed approach on an average 70% cases for both in short and long documents could successfully improve the ranking. In the future, we would like to explore different explanation units for the counterfactual setup.

## 7 Limitations

One of the limitations of this work is that we assume that top 10 or 20 words (based on tf-idf weights) within a document play the most important part in improving the rank of a document. However, theoretically speaking we should consider all the words present in a document to de-

604 termine the most influential words for a retrieval  
605 model. We have used top tf-idf words (Similar  
606 to statistical retrieval models) to reduce the com-  
607 putational complexity of our experiments and we  
608 have seen that increasing the number of top words  
609 doesn't affect the performance of the model that  
610 much.

## 611 8 Ethical Considerations

612 In this work, we have used publicly available  
613 search query log and document collection to  
614 demonstrate counterfactual explanation. No sen-  
615 sitive data was used in this experiment. As a result  
616 of this there is no particular ethical concern asso-  
617 ciated with this work. If there is any kind of bias  
618 present in the search log data that effect can be ob-  
619 served within our approach. However mitigating  
620 that bias was beyond the scope of this work

## 621 References

622 Avishek Anand, Lijun Lyu, Maximilian Idahl, Yumeng  
623 Wang, Jonas Wallat, and Zijian Zhang. 2022. Ex-  
624 plainable information retrieval: A survey.

625 Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng,  
626 Jianfeng Gao, Xiaodong Liu, Rangan Majumder,  
627 Andrew McNamara, Bhaskar Mitra, Tri Nguyen,  
628 Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Ti-  
629 wary, and Tong Wang. 2016. Ms marco: A human  
630 generated machine reading comprehension dataset.  
631 In *InCoCo@NIPS*.

632 Alexander Bondarenko, Maik Fröbe, Jan Heinrich  
633 Reimer, Benno Stein, Michael Völske, and Matthias  
634 Hagen. 2022. Axiomatic retrieval experimentation  
635 with ir axioms. In *Proc. of SIGIR 2022*, pages  
636 3131–3140.

637 Miguel Á Carreira-Perpiñán and Suryabhan Singh  
638 Hada. 2021. Counterfactual explanations for  
639 oblique decision trees: Exact, efficient algorithms.  
640 In *Proceedings of the AAAI conference on artificial  
641 intelligence*, volume 35, pages 6903–6911.

642 Nick Craswell, Bhaskar Mitra, Emine Yilmaz, and  
643 Daniel Campos. 2021. [Overview of the TREC 2020  
644 deep learning track](#). *CoRR*, abs/2102.07662.

645 Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel  
646 Campos, Jimmy Lin, Ellen M. Voorhees, and  
647 Ian Soboroff. 2023. [Overview of the trec 2022  
648 deep learning track](#). In *Text REtrieval Conference  
649 (TREC)*. NIST, TREC.

650 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
651 Kristina Toutanova. 2019. BERT: Pre-training of  
652 deep bidirectional transformers for language under-  
653 standing. In *NAACL-HLT*.

Gokhan Egri and Coskun Bayrak. 2014. [The role  
654 of search engine optimization on keeping the user  
655 on the site](#). *Procedia Computer Science*, 36:335–  
656 342. Complex Adaptive Systems Philadelphia, PA  
657 November 3-5, 2014. 658

Anett Erdmann, Ramón Arilla, and José M. Ponzoa.  
2022. [Search engine optimization: The long-term  
659 strategy of keyword choice](#). *Journal of Business Re-  
660 search*, 144:650–662. 661 662

Thibault Formal, Benjamin Piwowarski, and Stéphane  
Clinchant. 2021. [Splade: Sparse lexical and expansion  
663 model for first stage ranking](#). In *Proceedings  
664 of the 44th International ACM SIGIR Conference  
665 on Research and Development in Information Re-  
666 trieval*, page 2288–2292. 667 668

Maarten Grootendorst. 2020. [Keybert: Minimal key-  
669 word extraction with bert](#). 670

Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce  
Croft. 2016. A deep relevance matching model for  
ad-hoc retrieval. In *Proceedings of the 25th ACM  
International Conference on Information and  
Knowledge Management, CIKM '16*, page 55–64,  
New York, NY, USA. Association for Computing  
Machinery. 671 672 673 674 675 676 677

Jiafeng Guo, Yixing Fan, Xiang Ji, and Xueqi Cheng.  
2019. [Matchzoo: A learning, practicing, and develop-  
678 ing system for neural text matching](#). In *Proceed-  
679 ings of the 42nd International ACM SIGIR Confer-  
680 ence on Research and Development in Information  
681 Retrieval, SIGIR'19*, pages 1297–1300. 682 683

Faisal Hamman, Erfan Noorani, Saumitra Mishra,  
Daniele Magazzeni, and Sanghamitra Dutta. 2023.  
Robust counterfactual explanations for neural net-  
works with probabilistic guarantees. In *Proceed-  
ings of the 40th International Conference on Ma-  
chine Learning, ICML'23*. JMLR.org. 684 685 686 687 688 689

Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng,  
Alex Acero, and Larry Heck. 2013. Learning deep  
structured semantic models for web search using  
clickthrough data. In *Proceedings of the 22nd ACM  
International Conference on Information & Knowl-  
edge Management, CIKM '13*, page 2333–2338. 690 691 692 693 694 695

Kentaro Kanamori, Takuya Takagi, Ken Kobayashi,  
Yuichi Ike, Kento Uemura, and Hiroki Arimura.  
2021. [Ordered counterfactual explanation by  
696 mixed-integer linear optimization](#). In *Proceedings  
697 of the AAAI Conference on Artificial Intelligence*,  
698 volume 35, pages 11564–11574. 699 700 701

Amir-Hossein Karimi, Gilles Barthe, Borja Balle, and  
Isabel Valera. 2020. [Model-agnostic counterfactual  
702 explanations for consequential decisions](#). In *Inter-  
703 national Conference on Artificial Intelligence and  
704 Statistics*, pages 895–905. PMLR. 705 706

Omar Khattab and Matei Zaharia. 2020. [Colbert: Ef-  
707 ficient and effective passage search via contextual-  
708 ized late interaction over bert](#). In *Proceedings of the  
709*



**BM25:** BM25<sup>2</sup> is a statistical retrieval model where the similarity between a query and a document is computed based on the term frequency of the query words present in the document, document frequency of the query words and also the document length.

**DRMM:** Deep Relevance Matching Model (DRMM) Guo et al. (2016) is a neural retrieval model where the semantic similarity between each pair of tokens corresponding to a query and a document is computed to estimate the final relevance score of a document.

**DSSM:** Deep Semantic Similarity Model (DSSM) Huang et al. (2013) is another neural retrieval model which uses word hashing techniques to compute the semantic similarity between a query and a document.

**ColBERT:** Contextualized Late Interaction over BERT (ColBERT) (Khattab and Zaharia, 2020), is an advanced neural retrieval model which exploits late interaction techniques based on BERT (Devlin et al., 2019) based representations of both query and document for retrieval.

**MonoT5:** MonoT5 (Nogueira et al., 2020) is a sequence-to-sequence model fine-tuned to predict the relevance of a query-document pair.

**Splade:** Splade (Formal et al., 2021) (Sparse Lexical and Expansion Model for Information Retrieval) combines the sparse interpretability of traditional IR models (like BM25) with the semantic power of deep learning. Unlike dense retrieval models that rely on vector similarity in embedding space, SPLADE encodes queries and documents into sparse high-dimensional vectors—essentially performing learned term expansions in a way that mimics the inverted index structure used in classic IR systems.

## B Retrieval Performance of IR Models

We use Lin et al. (2021) toolkit for implementing BM25 and MonoT5 and Splade. For DRMM and DSSM, we use the implementation released by the study in Guo et al. (2019). For passage ranking we varied the parameters in a grid search and we took the configuration producing best MRR@10 value on TREC DL (Craswell et al., 2021) test set. For both DRMM and DSSM experiments on MSMARCO data, the parameters were set as suggested in (Wu et al., 2022b). The MRR@10 values are reported in Table 7 in Appendix B.

<sup>2</sup>[https://en.wikipedia.org/wiki/Okapi\\_BM25](https://en.wikipedia.org/wiki/Okapi_BM25)

For DRMM and DSSM, we use randomly chosen 100K query pairs from the MSMARCO training dataset to train the model.

The machine used to run counterfactual experiments on retrieval model has 1 A100 GPU and 40 GB memory.

Model	MRR@10	
	MSMARCO Passage	MSMARCO Document
BM25	0.1874	0.2184
DRMM	0.1623	0.1168
DSSM	0.1320	0.1168
ColBERT	0.3481	0.3469
MonoT5	0.3904	0.3827
Splade	0.3813	0.3721

Table 2: Retrieval Model Performance on MSMARCO passage and document

## C Dataset Statistics

The dataset statistics for all the experiments are given in Table 3

		MS MARCO Passage	MS MARCO Document	TREC Robust
Query	Avg Length	5.9	6.9	7.18
Document	Avg Length	64.9	1134.2	150.12
Query	#Instances	100	100	100
Document	#Instances	500	500	500

Table 3: Dataset Details for Counterfactual Setup

## D Example of Input and Output to Classifier

Given an input query, we employ a Lucene-Searcher with MSMARCO Index to retrieve the Top-K documents. The feature vector construction process follows these steps:

For each document, we:

1. Extract the top n words based on their Imp(d) values
2. Construct a vocabulary  $V$  as the union of all top 10 words across documents
3. Note that  $|V|$  typically falls in the range of 150-180 words

The feature vector for each document has dimension  $|V|$ , where each component represents the value from the Imp(d) of the corresponding word from the vocabulary. Formally:

$$d_{vec} \in R^{|V|}$$

Labels are assigned according to the following criterion:

$$\text{label} = \begin{cases} 1 & \text{for top } K \text{ documents} \\ 0 & \text{for remaining documents} \end{cases}$$

Example feature vectors and their corresponding counterfactuals generated using (Mothilal et al., 2020) are shown in Table 5. Since  $|V|$  is 150 in our experiments, hence in Table 5 we have only shown the term frequencies of the words present in each document. For other words the terms frequency values will be zero in  $d_{vec}$ .

Existing Explanation Methods	Word Overlap
PointWise Explanation (Singh and Anand, 2019)	21.46%
ListWise Explanation (Lyu and Anand, 2023)	9.57%

Table 4: Comparison of CFIR with Existing ExIR Approaches

## E Scalability Issues

There can be concerns regarding the feasibility of training a classifier per document. To address this, we propose and evaluate an alternative and more efficient strategy in which a single classifier is trained per query, rather than per document. Concretely, for a given query, we train one classifier using: (i) all documents for which explanations are required (let this number be  $x$ ); (ii) their nearest neighbors, which contribute to the non-relevant document set; and (iii) the top-k retrieved documents. To balance the number of relevant and non-relevant training instances, we construct the non-relevant set with a total size of  $2k$ , where we sample  $2k/x$  nearest neighbors from each document for which explanations are generated.

The results of this per-query training strategy are reported in Table 6. As shown, this substantially faster approach achieves performance comparable to that reported in Table 2, where a separate classifier was trained for each document. This demonstrates that our method remains effective while significantly reducing the computational overhead, directly addressing the reviewers’ feasibility concerns.

## F Existing EXIR approaches vs. CFIR

The existing literature aims to explain the significance of a document, a set of documents, or a pair of documents through various explanation methods. Nonetheless, our proposed approach diverges fundamentally from prior work in that we seek to demonstrate how the absence or frequency of certain tokens impacts document relevance. In this section, we examine whether there is any intersection between the two sets of tokens described earlier.

**Pointwise Explanation Approach** As outlined in Section 2, existing pointwise explanation methods elucidate why a specific document aligns with a given query within a retrieval model. Similarly, our proposed approach operates on individual documents and queries, albeit with a distinct objective. Here, we analyze the overlap between the explanations generated by the pointwise explanation method and those derived from our model, as presented in Table 7. This comparison was conducted on 50 pairs of documents.

**Listwise Explanation Approach** In Section 2, it is explained that listwise explanations typically aim to demonstrate the relevance of a list of documents to a given query. In listwise setup, one set of explanation terms are extracted for a list of documents, a query, and a retrieval model. Conversely, in our approach, we generate distinct explanations for each query-word pair. Therefore, to compare listwise explanations with our method, we aggregate all individual explanations obtained for each document-query pair in the list to create a unified explanation set for the entire list corresponding to a query. The resulting overlap is presented in Table 7.

## G Counterfactual Optimization Framework

The different parts of Equation 2 are described here. The  $y_{loss}$  in Equation 2 is a hinge loss function as defined in Equation 4. In Equation 4  $z$  is  $-1$  when  $y = 0$  otherwise,  $z = 1$ .  $logit(f_{\{M,q\}}(c_i))$  is the logit values obtained from the classifier ( $f_{\{M,q\}}$ ) when the counterfactual  $c_i$  is given as input.

$$y_{loss} = \max(0, 1 - z * \logit(f_{\{M,q\}}(c_i))) \quad (4)$$

The distance function ( $dist(c_i, d)$ ) in Equation 2 is computed using the formula given in Equation 5. In Equation 5,  $V$  represents the vocabulary set used to represent the document vectors ( $d_{vec}$ ). In Equation 5, the value of  $I$  is equal to 1 if the corresponding term is present in both the counterfactual input  $c$  and the original input  $d$ , otherwise it is set to 0.

$$dist(c, d) = \sum_{p=1}^V I(c_p \neq d_p) \quad (5)$$

The diversity in above equation is defined by the formula described in Equation 6. In equation 6,

docID	Feature Vector
3686955	[prohibition:2.0, amendment:2.0, under:1.0, dwindled:1.0, eighteenth:1.0, repeal:1.0, repealed:3.0, states:1.0, 1933: 1.0, ratification: 1.0]
6159679	[membrane:5.0, lipids:3.0, remainder:2.0, proteins:3.0, biochemical:2.0, 80:2.0, role:2.0, percent:2.0]
5217641	[waves:6.0, transverse:5.0, electromagnetic:3.0, oscillations:2.0, vibrations:2.0, travel:2.0, radiation:2.0, angles:2.0, transfer:2.0, types:3.0]

Table 5: Sample Feature Vector Corresponding to three different documents

Model Description		MS MARCO Passage			MS MARCO Document			Trec Robust		
Retrieval Model	Classifier	FD(%)	Avg. New Words	Avg. Query Overlap	FD(%)	Avg. New Words	Avg. Query Overlap	FD(%)	Avg. New Words	Avg. Query Overlap
<i>CFIR<sub>MonoT5</sub></i>	LR	<b>81.16%</b>	12.45	63%	<b>73%</b>	<b>13.13</b>	63%	<b>74%</b>	<b>10.45</b>	67%
<i>CFIR<sub>Splade</sub></i>	RF	78%	11.23	62%	69%	12.11	60%	71%	9.81	65%
<i>CFIR<sub>Splade</sub></i>	LR	<b>76.92%</b>	12.15	63.4%	<b>68%</b>	<b>11.33</b>	64%	<b>70.11%</b>	<b>8.91</b>	67%

Table 6: CFIR model Performance for MonoT5 and Splade in MSMARCO Passage and Document Collection and TREC Robust. The Best Performing Counterfactual Explanation Method for every retrieval model is boldfaced; the overall best performance across all rows is underlined. All the results reported in Table 1 are statistically significant with  $p < 0.05$ .

Existing Explanation Methods	Word Overlap
PointWise Explanation (Singh and Anand, 2019)	21.46%
ListWise Explanation (Lyu and Anand, 2023)	9.57%

Table 7: Comparison of CFIR with Existing ExIR Approaches

$K_{i,j}$  is equal to  $\frac{1}{1+dist(c_i, c_j)}$ .  $dist(c_i, c_j)$  calculates the distance between two counterfactuals  $c_i$  and  $c_j$ .

$$div(c_1, \dots, c_k) = \sum_{i,j} det(K_{i,j}) \quad (6)$$

## H Parameters for Counterfactual Setup

The value of  $\lambda_1$  and  $\lambda_2$  is set to 1 and 0.5 respectively in Equation 2. The value of  $k$  in Equation 2 is set to  $k = 3$ . In all our experiments in Table 1, we have observed that for  $K = 3$  and onward we have always found a counterfactual explanation for each query-document pair where only words were added for the desired counterfactual outcome.

## I Example of Counterfactuals Produced by CFIR Setup

The words shown in Table 8 have improved the ranking of a docID with respect to the queries shown.

## J Adversarial Attacks vs. Counterfactual Explanation

Here we show the performance of our proposed counterfactual explanation approach with an existing adversarial model named PRADA (Wu et al., 2022a). We use the MSMARCO passage dataset as the target corpus. We use same test set (as described in Table 3) as used in the first column of

Table 1 in this experiment. Table 9 shows the results in terms of Fidelity score.

## K Implementation of Imp(d)

We explored three ways to compute the top  $n$  words from each document. Each one of them is described as follows.

**TF-IDF Approach:** In this approach we choose top  $n$  words from a document based on their TF-IDF weight.

**KEYBERT Approach:** In this approach we use the model proposed in (Grootendorst, 2020) to extract keywords from a string.

**BERT-Based Similarity(BERTSim):** In this approach we compute the similarity between the BERT based representation of the query text and each token of the document and then we sort all the tokens based on the similarity.

Table 10 shows the performance of the above-mentioned three approaches in MSMARCO passage dataset and ColBERT retrieval model.  $n = 10$  for the experiments shown in Table 10. From Table 10, we can conclude that the BERT-based similarity approach works the best for the  $Imp(d)$  function. hence for all the results reported in Table 1, we use the BERTSim approach in the  $Imp(d)$  function.

## L User Study

In the user study we didn't record any personal information of any of the users. We only recorded their judgment about the output of the proposed methodology for the study. We have also used data which is publicly available for IR research. Hence no ethics approval was required for the study. All

Retrieval Model	Query Text	docId	Explanation Terms
DRMM	What law repealed prohibition ?	3686955	working, strict, Maine, 1929, law, resentment, New York City, Irish, immigrant, prohibition, repeal, fall, Portland, temperance, riot, visit
DSSM	What is the role of lipid in the cell?	6159679	phospholipid, fluidity, storage, triglyceride, fatty receptor
ColBERT	what type of wave is electromagnetic?	5217641	directly ,oscillations, medium, wave, properties, speed
MonoT5	what is a caret?	6338711	display, diamond, weight
Splade	which vitamins help heal bruises?	3465680	minerals, body, eat, cut

Table 8: CFIR explanation terms for DRMM, DSSM, ColBERT, MonoT5 and Splade in MS MARCO passage.

Retrieval Model	FD in PRADA	FD in CFIR
ColBERT	74%	75%
MonoT5	80%	82%

Table 9: Performance of CFIR vs. Adversarial Attack Model PRADA (Wu et al., 2022a)

$Imp(d)$ Approach	FD
TFIDF	74%
KeyBERT	70%
BERTSim	<b>75%</b>

Table 10: Performance of Different Approaches in  $Imp(d)$ .

the researchers were made aware of the of the use of their assessment in this research.

1045  
1046