

---

# A Neural Corpus Indexer for Document Retrieval

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Current state-of-the-art document retrieval solutions mainly follow an index-  
2 retrieve paradigm, where the index is hard to be optimized for the final retrieval  
3 target. In this paper, we aim to show that an end-to-end deep neural network unifying  
4 training and indexing stages can significantly improve the recall performance  
5 of traditional methods. To this end, we propose Neural Corpus Indexer (NCI), a  
6 sequence-to-sequence network that generates relevant document identifiers directly  
7 for a designated query. To optimize the recall performance of NCI, we invent a  
8 prefix-aware weight-adaptive decoder architecture, and leverage tailored techniques  
9 including query generation, semantic document identifiers and consistency-based  
10 regularization. Empirical studies demonstrated the superiority of NCI on a com-  
11 monly used academic benchmark, achieving +51.9% relative improvement on  
12 NQ320k dataset compared to the best baseline.

## 13 1 Introduction

14 Document retrieval and ranking are two key stages for a standard web search engine [46, 27]. First,  
15 the document retrieval stage retrieves candidate documents relevant to the query, and then, the ranking  
16 stage gives a more precise ranking score for each document. The ranking stage is often fulfilled  
17 by a deep neural network, taking each pair of query and document as input and predicting their  
18 relevance score. Nevertheless, a precise ranking model is very costly, while typically only a hundred  
19 or thousand candidates per query are affordable in an online system. Therefore, the recall performance  
20 of document retrieval stage is very crucial to the effectiveness of web search engine.

21 Existing document retrieval methods can be divided into two categories, namely *term-based* and  
22 *semantic-based* approaches [18]. Term-based retrieval approaches [8, 48] build an inverted index  
23 for the entire web corpus, but they hardly capture document semantics and fail to retrieve similar  
24 documents in different wordings. Thus, semantic-based approaches [46, 29] are proposed to alleviate  
25 this discrepancy. First, they learn dense representations for both queries and documents through a  
26 twin-tower architecture; then Approximate Nearest Neighbor (ANN) search is applied to retrieve  
27 relevant documents for the designated query. Despite of their success in real applications, these  
28 approaches can not fully leverage the power of deep neural networks for the following reasons.  
29 First, a single embedding vector has limited capacity to memorize all semantics in a document,  
30 and it performs even worse than term-based methods in applications that heavily rely on exact  
31 match [30]. Second, the model is unable to incorporate deep query-document interactions. Because  
32 ANN algorithms theoretically require a strong assumption for the Euclidean space, we have to adopt  
33 simple functions such as cosine similarity to capture the query-document interactions [16].

34 Given the above limitations, several research works have explored end-to-end models that directly  
35 retrieve relevant candidates without using an explicit index. Gao et al. [16] proposed a Deep Retrieval  
36 (DR) framework for item recommendation, which learned a retrievable structure with historical  
37 user-item interactions. Nevertheless, it is more challenging to design a universal model for semantic  
38 text retrieval, as we need to leverage the power of both pre-trained language models and deep retrieval

39 networks simultaneously. Tay et al. [41] proposed Differentiable Search Index (DSI), a text-to-text  
40 model that maps queries directly to relevant docids. To the best of our knowledge, this is the first  
41 attempt to propose a differentiable index for semantic search. However, [the vanilla transformer  
42 decoder in DSI does not fully leverage the hierarchical structures of document identifiers](#), and the  
43 model is pruned to over-fitting with limited training data. Furthermore, Bevilacqua et al. [4] proposed  
44 SEAL by leveraging all n-grams in a passage as its identifiers. But for long documents, it is hard to  
45 enumerate all possible n-grams. In general, the recall performance of end-to-end document retrieval  
46 remains a large room to be improved.

47 In this paper, we show that the traditional text retrieval frameworks can be fundamentally changed  
48 by a unified deep neural network with tailored designs. To this end, we propose a Neural Corpus  
49 Indexer (NCI), which supports end-to-end document retrieval by a sequence-to-sequence neural  
50 network. The model takes user query as input, generates query embedding through the encoder, and  
51 outputs the identifiers of relevant documents using the decoder. It can be trained by ground-truth  
52 and augmented query-document pairs. During inference, the top  $N$  documents are retrieved directly  
53 via beam search based on the decoder. Designing and training such a model is non-trivial, so we  
54 propose several crucial techniques to ensure its effectiveness. First, to get sufficient query-document  
55 pairs for training, we leverage a query generation network to obtain possible pairs of queries and  
56 documents. Second, we utilize the hierarchical  $k$ -means algorithm to generate a semantic identifier  
57 for each document. Third, we design a prefix-aware weight-adaptive decoder to replace the vanilla  
58 one in a sequence-to-sequence architecture. Specifically, the same token will be assigned different  
59 embedding vectors at different positions in the identifiers, while another transformer-based adaptive  
60 module is applied to the classification weights for token prediction in the context of a certain prefix.  
61 This makes the classifiers customized to different prefixes when decoding along the hierarchical tree  
62 structure. Besides, a consistency-based regularization loss is taken for training both encoder and  
63 decoder networks to mitigate the over-fitting problem.

64 Our NCI design solves the limitations of traditional index-retrieve pipelines from multiple perspec-  
65 tives. On one hand, a whole neural network model replaces the traditional inverted index or vector  
66 search solutions. It can be optimized end-to-end using realistic query-document pairs, which fully  
67 capture both term-based and semantic-based features and is adaptive to the changing of workloads.  
68 On the other hand, the model is able to capture deep interactions between queries and documents via  
69 encoder-decoder attention, which enlarges the capacity of vector-based representations. Moreover,  
70 NCI achieves much better ranking results than ANN-based approaches as it is optimized directly by  
71 the final target. Thus, it can be served as an end-to-end retrieval solution and release the burden of  
72 re-ranking for a long candidate list.

73 In addition to the superior performance, the invention of Neural Corpus Indexer is also promising  
74 from the perspective of system design. As nowadays, ranking and query-answering modules are  
75 already implemented by neural networks, NCI finishes the last piece of puzzle for the next-generation  
76 information retrieval system based on a unified differentiable model architecture. This reduces the  
77 dependency among different sub-modules, while the process of system deployment and maintenance  
78 could be greatly eased.

79 Our **contributions** are highlighted as follows.

- 80 • For the first time, we demonstrate that an end-to-end differentiable document retrieval model  
81 can significantly outperform both inverted index and dense retrieval solutions. This finding will  
82 inspire research on further steps towards the next-generation search systems, for instance, unifying  
83 informational retrieval, ranking and question answering in a single differentiable framework.
- 84 • We design a sequence-to-sequence model, named Neural Corpus Indexer (NCI), which generates  
85 relevant document identifiers directly for a specific query. In our experiments, the proposed NCI  
86 model improves the state-of-the-art performance of existing methods by a significant margin,  
87 achieving +51.9% and +19.2% relative enhancement for Recall@1 and Recall@10 respectively  
88 on NQ320k dataset. Also, NCI itself can achieve a competitive MRR score without using an  
89 explicit ranking model.
- 90 • We propose a novel decoder architecture, namely *prefix-aware weight-adaptive (PAWA)* decoder,  
91 to generate document identifiers. As verified by ablation studies, this invention is very crucial for  
92 NCI to achieve an outstanding performance. Moreover, query generation, semantic document  
93 identifiers and consistency-based regularization are all accountable for the superior capability of  
94 Neural Corpus Indexer.

## 95 2 Related work

96 In this section, we briefly introduce the related works and leave more discussions about the traditional  
97 web search techniques in the Appendix A.

98 **Sparse retrieval.** Traditional document retrieval methods are based on *Sparse Retrieval*, which is  
99 built upon inverted index with term matching metrics such as TF-IDF [37], query likelihood [26]  
100 or BM25 [36]. In industry-scale web search, BM25 is a difficult-to-beat baseline owing to its  
101 outstanding trade-off between accuracy and efficiency. In recent years, there are some attempts  
102 to incorporate the power of neural networks into inverted index. The Standalone Neural Ranking  
103 Model (SNRM) [47] learns high-dimensional sparse representations for query and documents, which  
104 enables the construction of inverted index for efficient document retrieval. Doc2Query [33] predicts  
105 relevant queries to augment the content of each document before building the BM25 index, and  
106 DocT5Query [32] improves the performance of query generation by the pre-trained language model  
107 T5 [5]. Furthermore, DeepCT [8] calculates context-aware term importance through neural networks  
108 to improve the term matching metrics of BM25.

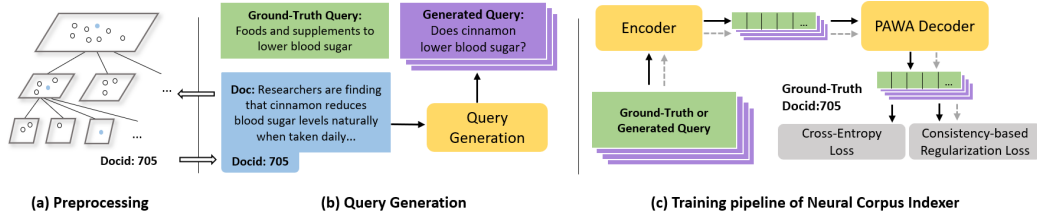
109 **Dense retrieval.** Another line of research lies in *Dense Retrieval*, which presents query and documents  
110 in dense vectors and models their similarities with inner product or cosine similarity. These methods  
111 benefit from recent progresses of pre-trained language models, such as BERT [13] and RoBERTa [28]  
112 to obtain dense representations for queries and documents. At inference time, efficient Approximate  
113 Nearest Neighbor (ANN) search algorithms, such as k-dimensional trees [3], locality-sensitive  
114 hashing [9], and graph-based indexes (e.g., HNSW [31], DiskANN [23] and SPANN [7]) can be  
115 utilized to retrieve relevant documents within a sublinear time. Besides, Luan et al. [30] analyze the  
116 limited capacity of dual encoders, and propose a combination of sparse and dense retrieval methods  
117 with multi-vector encoding to achieve better search quality.

118 **Autoregressive retrieval.** The other way to approach retrieval is utilizing an end-to-end autoregressive  
119 models. Firstly, several efforts have been done on entity linking [12, 11, 10], which can be regard  
120 as a special type of retrieval task, e.g., using an entity to ask the posed question. Recently, different  
121 from the entity linking task, Tay et al. [41] proposed the DSI (differentiable search index) to generate  
122 relevant document identifiers directly according to the query. Bevilacqua et al. [4] employ the  
123 autoregressive model to generate the relevant words for a query and utilize the generated string to  
124 retrieve relevant documents. Besides, the Deep Retrieval (DR) [16] approach for recommendation is  
125 also related to this category, which learns a deep retrievable network with user-item clicks and gets  
126 rid of the ANN algorithms based on the Euclidean space assumption.

127 **Pre-trained language models.** Recently, pre-trained Language Models (LMs), such as BERT [13]  
128 and RoBERTa [28], have led to a revolution in web search techniques. The representation vectors  
129 for all documents can be calculated and indexed offline. In the online serving stage, it calculates the  
130 representation vector for the input query and applies a crossing layer to calculate the relevance score  
131 between each query and document. The crossing layer usually adopts simple operators such as cosine  
132 similarity or a single feed-forward layer to retain a high efficiency. Gao et al. [14] find that a standard  
133 LMs’ internal attention structure is not ready-to-use for dense encoders and propose the Condenser  
134 to improve the performance of dense retrieval. Moreover, ANCE [45] leverages hard negatives to  
135 improve the effectiveness of contrastive learning, which generates better text representations for the  
136 retrieval tasks.

## 137 3 Neural corpus indexer

138 The neural corpus indexer (NCI) is a sequence-to-sequence neural network model. The model takes  
139 *query* as input and outputs the most relevant *document identifier (docid)*, which can be trained by a  
140 large collection of  $\langle query, docid \rangle$  pairs. The documents are encoded into semantic *docids* by the  
141 hierarchical  $k$ -means algorithm [19], which makes similar documents have “close” identifiers in the  
142 hierarchical tree. As shown in Figure 1, NCI is composed of three components, including *Query*  
143 *Generation*, *Encoder* and *Prefix-Aware Weight-Adaptive (PAWA) Decoder*. Query generation is imple-  
144 mented by a sequence-to-sequence transformer model [43] that takes as an input the document terms  
145 and produces a query as output [33]. The encoder, following the standard transformer architecture, is  
146 composed of  $N_1$  stacked transformer blocks, which outputs the representation for an input query. For  
147 the decoder network, we stack  $N_2$  transformer layers. To better align with the hierarchical nature of  
148 the semantic identifiers, we propose a weight adaptation mechanism based on another transformer



**Figure 1:** Overview of Neural Corpus Indexer (NCI). (a) Preprocessing. Each document is represented by a semantic identifier via hierarchical  $k$ -means. (b) Query Generation. Queries are generated for each document based on the content. (c) The training pipeline of NCI. The model is trained over augmented  $\langle \text{query}, \text{docid} \rangle$  pairs through a standard transformer encoder and the proposed Prefix-Aware Weight-Adaptive (PAWA) Decoder.

149 to make the decoder aware of semantic prefixes. At inference time, the top  $N$  relevant documents  
 150 can be easily obtained via beam search. Due to the hierarchical property of semantic identifiers, it is  
 151 relatively easy to constrain the beam search decoding process on the prefix tree so that only valid  
 152 identifiers will be generated.

### 153 3.1 Representing document with semantic identifiers

154 NCI generates document identifiers solely based on the input query without explicit document content,  
 155 which is difficult when the size of the corpus is very large. Thus, we aim to inject useful priors into  
 156 the identifiers, so that the semantic information of documents can be considered in the tree-based  
 157 decoding process. In other words, we hope the documents with similar information have close  
 158 *docids* to facilitate the learning process of NCI. To achieve this, we leverage the hierarchical  $k$ -means  
 159 algorithm to encode documents. As shown in Figure 1(a), given a collection of documents be indexed,  
 160 all documents are first classified into  $k$  clusters by using their representations encoded by BERT [13].  
 161 For cluster with more than  $c$  documents, the  $k$ -means algorithm is applied recursively. For each  
 162 cluster containing  $c$  documents or less, each document is assigned a number starting from 0 to at  
 163 most  $c-1$ . In this way, we organize all documents into a tree structure  $T$  with root  $r_0$ . Each document  
 164 is associated with one leaf node with a deterministic routing path  $l = \{r_0, r_1, \dots, r_m\}$  from the root,  
 165 where  $r_i \in [0, k)$  represents the internal cluster index for level  $i$ , and  $r_m \in [0, c)$  is the leaf node.  
 166 The semantic identifier for a document is concatenated by the node indices along the path from  
 167 root to its corresponding leaf node. For documents with similar semantics, the prefixes of their  
 168 corresponding identifiers are likely to be the same. For simplicity, we set  $k = 10$  and  $c = 10$  in  
 169 all experiments, leaving the optimization of these hyper-parameters to future work. The detailed  
 170 procedure of hierarchical  $k$ -means will be described in Algorithm 1 in the Appendix B.2.

### 171 3.2 Query generation

172 One challenge of generating document identifiers by single query input is how to make the identifiers  
 173 aware of the document semantics. Since the content of each document is not explicitly known at  
 174 inference, it must be incorporated into the model parameters during training. To facilitate the training  
 175 process, we generate a bunch of queries with a query generation module and bind the information  
 176 of document content through training the sequence-to-sequence model with generated queries and  
 177 their document identifiers. We adopt a standard sequence-to-sequence transformer [43] based on the  
 178 implementation of Doc2Query [1], which takes as an input the document terms and produces relevant  
 179 queries via random sampling. Note that we use random sampling instead of beam search to ensure  
 180 the diversity of generated queries.

### 181 3.3 Prefix-aware weight-adaptive decoder

182 The probability of generating a document identifier can be written as follows:

$$p(l|x, \theta) = \prod_{i=1}^m p(r_i|x, r_1, r_2, \dots, r_{i-1}, \theta_i), \quad (1)$$

183 where  $r_i$  is the  $i$ -th token in the current identifier;  $x$  is the representation output from encoder;  $\theta$   
 184 denotes the total parameters and  $\theta_i$  is the parameter for the  $i$ -th step.

185 This probability can be modeled by a transformer-based  
 186 decoder. For an internal node with level  $i$ , the probabili-  
 187 ty is calculated by:

$$h_i = \text{TransformerDecoder}(x, h_1, h_2, \dots, h_{i-1}; \theta_i), \quad (2)$$

$$188 \quad p(r_i|x, r_1, r_2, \dots, r_{i-1}, \theta_i) = \text{Softmax}(h_i W). \quad (3)$$

189 Here  $h_i$  is the hidden representation for step  $i$ , which  
 190 is calculated by a multi-head attention over encoder  
 191 representation  $x$  and token representations of previous  
 192 decoding steps. The linear classification weight is de-  
 193 noted by  $W \in \mathbb{R}^{d \times v}$ ,  $d$  is the hidden dimension size  
 194 and  $v$  is the vocabulary size of identifiers.

195 As the encoder and decoder utilize distinct vocabulary spaces, we do not share the embedding space  
 196 for their tokens. Different from a standard decoding task, the meanings of the same token appearing  
 197 at different places of the same identifier are different, as they correspond to different clusters in the  
 198 hierarchical tree structure. For instance, the “5<sub>2</sub>” and “5<sub>3</sub>” of the same identifier “3<sub>1</sub>5<sub>2</sub>5<sub>3</sub>” correspond  
 199 to different semantic meanings. Moreover, the same token in the same position may have different  
 200 semantics with different prefixes. For example, in identifiers “1<sub>1</sub>1<sub>2</sub>5<sub>3</sub>” and “2<sub>1</sub>4<sub>2</sub>5<sub>3</sub>”, the same  
 201 token “5<sub>3</sub>” has different semantics in two different identifiers, as they are routed from different prefix  
 202 paths. The two properties of the hierarchical semantic identifiers motivate us to design the novel  
 203 **Prefix-Aware Weight-Adaptor (PAWA)** decoder.

204 Unlike a standard transformer decoder, the probabilities at different tree levels, such as  
 205  $p(r_i|x, r_{1..i-1}, \theta_i)$  and  $p(r_j|x, r_{1..j-1}, \theta_j)$  when  $i \neq j$ , do not share parameters with each other.  
 206 To distinguish different semantic levels, we concatenate the position and token values as input for  
 207 each decoding step, as shown in the left corner of Figure 2. Specifically, we have “(1, 3)(2, 5)(3, 5)”  
 208 for the semantic identifier “3<sub>1</sub>5<sub>2</sub>5<sub>3</sub>”, while “(2, 5)” and “(3, 5)” represent different tokens in the  
 209 vocabulary space. As the token embedding and linear classification layers share the same weights, the  
 210 same token value in different positions would correspond to different model parameters. Moreover,  
 211 to reflect the influence of different prefixes, we expect the linear classification layer to be aware of  
 212 different prefixes for predicting a specific token. Concretely, instead of using the same projection  
 213 weight  $W$  in the linear classification layer, we employ the prefix-aware adaptive weights for each  
 214 token classifier, which can be calculated by another transformer decoder,

$$W_{ada}^i = \text{AdaptiveDecoder}(e; r_1, r_2, \dots, r_{i-1}) W_i \quad (4)$$

215 where  $e$  is the query embedding vector taken as initial input to the transformer decoder;  $\{r_t|t \in$   
 216  $(1, 2, \dots, i-1)\}$  are prefix tokens before the  $i$ -th position, AdaptiveDecoder stacks  $N_3$  transformer de-  
 217 coding layers with dimension  $d$ , and  $W_{ada}^i \in \mathbb{R}^{d \times v}$  is the adapted weight matrix for the corresponding  
 218 classifier. Finally, the  $i$ -th token in the given prefix can be predicted by  $\text{Softmax}(h_i W_{ada}^i)$ .

219 For instance, to predict the third tokens in the identifiers “(1,3)(2,1)(3,5)” and “(1,2)(2,4)(3,5)” respec-  
 220 tively, the corresponding adaptive weights are derived separately for different prefixes “(1,3)(2,1)” and  
 221 “(1,2)(2,4)”. As we already know the previous tokens for each position in the teacher forcing setting,  
 222 the prefix-aware adaptive weights can be calculated and trained in parallel in different positions while  
 223 adding little burden to the entire model.

### 224 3.4 Training and inference

225 **Consistency-based regularization.** To alleviate over-fitting, we employ a consistency-based reg-  
 226 ularization loss for training each decoding step. Given an input query  $q$ , we denote the model  
 227 probabilities predicted by **two forward passes with independent dropouts** as  $p_1(r_i|E(q), r_1, \dots, r_{i-1}, \theta_i)$   
 228 and  $p_2(r_i|E(q), r_1, \dots, r_{i-1}, \theta_i)$  respectively, where  $E(\cdot)$  denotes the encoder network. The consistency-  
 229 based regularization loss tries to regularize the model prediction by minimizing the bidirectional  
 230 Kullback-Leibler (KL) Divergence between two output probabilities with random dropout. The  
 231 regularization loss of query  $q$  for the  $i$ -th decoding step is defined as,

$$\mathcal{L}_{reg} = \frac{1}{2} \left( \sum_{i=1}^m \mathcal{D}_{KL}(p_1(r_i|E(q), r_1, \dots, r_{i-1}, \theta_i) \| p_2(r_i|E(q), r_1, \dots, r_{i-1}, \theta_i)) \right. \\ \left. + \mathcal{D}_{KL}(p_2(r_i|E(q), r_1, \dots, r_{i-1}, \theta_i) \| p_1(r_i|E(q), r_1, \dots, r_{i-1}, \theta_i)) \right). \quad (5)$$

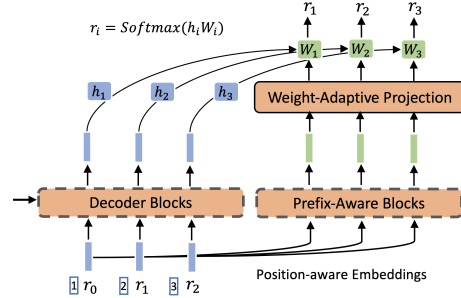


Figure 2: Overview of the Prefix-Aware Weight-Adaptive (PAWA) Decoder.

232 **Training loss.** Given a set of training examples  $\mathcal{D} = \{(q, d)\}$  composed of queries (training queries  
233 and augmented queries) and document identifiers, the loss function can be written as follows:

$$\mathcal{L}(\theta) = \sum_{(q,d) \in \mathcal{D}} (\log p(d|E(q), \theta) + \alpha \mathcal{L}_{reg}), \quad (6)$$

234 where  $p(d|E(q), \theta)$  denotes the probability of generating  $d$  with  $q$  as the input. The first part is  
235 the seq2seq cross-entropy loss with teacher forcing and the second part is the consistency-based  
236 regularization loss summed by all decoding steps. The whole process formulates a sequence-to-  
237 sequence neural network, which can be optimized end-to-end via gradient descent. The hyper-  
238 parameter  $\alpha$  denotes a scaling factor of regularization loss, which will be analyzed in Section 4.4.

239 **Inference via beam search.** In the inference stage, we calculate the query embedding through the  
240 encoder network and then perform the beam search on the decoder network. Due to the hierarchical  
241 nature of *docid*, it is convincing to constrain the beam search decoding process with a prefix tree,  
242 which in turn only generates the valid identifiers. The time complexity of beam search is  $O(LBF)$ ,  
243 where  $L$  is the max length of identifiers (the depth of tree),  $B$  is the beam size and  $F$  is the max  
244 fanout of the tree (10 in our experiments). Given a balanced tree structure built by a corpus with  
245  $M$  documents, the average time complexity for beam search is  $O(B \log M)$ . We leave detailed  
246 descriptions of the constrained beam search algorithm in Appendix B.3.

## 247 4 Experiments

248 In this section, we empirically verify the performance of NCI and the effectiveness of each com-  
249 ponent on the document retrieval task, which generates a ranking list of documents in response to  
250 a query. In the following, we discuss the datasets and evaluation protocol in Section 4.1, describe  
251 the implementation details and baseline methods in Section 4.2, and present empirical results and  
252 analyses in Section 4.3 and 4.4 respectively.

### 253 4.1 Datasets & evaluation metrics

254 **Datasets.** Following DSI [41] and SEAL [4], we conduct our experiments on the Natural Questions  
255 [25] dataset. *Natural Questions* (NQ) dataset was introduced by Google in 2019 [25]. The version  
256 we use is often referred to as NQ320k, which consists of 320k query-document pairs, where the  
257 documents are gathered from Wikipedia pages and the queries are natural language questions. We  
258 use its predetermined training and validation split for evaluation.

259 **Metrics.** We use widely accepted metrics for information retrieval, including Recall@ $N$  and Mean  
260 Reciprocal Rank (MRR). Recall@ $N$  measures how often the desired document is hit by the top- $N$   
261 retrieved candidates. MRR calculates the reciprocal of the rank at which the first relevant document is  
262 retrieved. A high recall means that the ground truth document is contained in the retrieved candidate  
263 list, while a high MRR indicates that the corresponding document has already been ranked at the top  
264 position without a need for re-ranking.

### 265 4.2 Implementation details

266 **Hierarchical semantic identifier.** For semantic identifiers, we apply hierarchical  $k$ -means algorithm  
267 over the document embeddings obtained through a 12 layers BERT model with pre-trained parameters  
268 provided by the HuggingFace [44]. For each hierarchical layer, we employ the default  $k$ -means  
269 algorithm implemented in scikit-learn [34] with  $k = 10$ . For simplicity, the recursion terminal  
270 condition is also set as  $c = 10$ .

271 **Query generation.** We leverage the pre-trained model docT5query [32] for query generation. We  
272 provide all documents in the NQ320k dataset to predict augmented query-document pairs. For each  
273 document, we generate 10 queries with the first 512 input tokens of the document as the input and  
274 constrain the maximum length of the generated query as 64.

275 **Training and inference.** Neural Corpus Indexer (our approach) are implemented with python 3.6.10,  
276 PyTorch 1.8.1 and HuggingFace transformers 3.4.0. We utilize the parameters of the T5 pre-trained  
277 model [5] to initialize the encoder and randomly initialize the PAWA decoder. All NCI experiments  
278 are based on a learning rate  $2 \times 10^{-4}$  for encoder and  $1 \times 10^{-4}$  for decoder with a batch size 16 per  
279 GPU. We set the scaling factor of the consistency-based regularization loss as  $\alpha = 0.015$ , and the  
280 dropout ratio is 0.1. For inference, we apply the partial beam search algorithm to the trained seq2seq

**Table 1:** Performance comparison on NQ320k retrieval task

Method	Recall@1	Recall@10	Recall@100	MRR@100
<b>Neural Corpus Indexer (Ours)</b>	<b>88.72</b>	<b>95.84</b>	<b>97.43</b>	<b>91.59</b>
DSI (T5-Base)	27.40	56.60	–	–
DSI (T5-XXL)	40.40	70.30	–	–
SEAL (BART-Base)	26.55	53.61	72.67	35.64
ANCE (FirstP)	51.33	80.33	<b>91.78</b>	61.71
ANCE (MaxP)	52.63	<b>80.38</b>	91.31	62.84
BERT + BruteForce	28.65	53.42	73.16	36.60
BERT + ANN (Faiss)	27.92	53.63	73.01	37.08
BM25 + DocT5Query	<b>58.39</b>	75.76	89.51	<b>64.53</b>
BM25	30.23	47.02	68.54	36.26

281 model. We set the length penalty and the beam size as 0.3 and 100 respectively. All experiments are  
 282 based on a cluster of NVIDIA V100 GPUs with 32GB memory. Each job takes 8 GPUs, resulting in  
 283 a total batch size of 128 ( $16 \times 8$ ).

284 **Baselines.** We evaluate BM25 on both raw documents and those augmented by DocT5Query. The  
 285 performance of DSI [40] is referred from its original paper as the implementation has not been  
 286 open-sourced. To avoid the difference in data processing, we reproduce SEAL [4] and ANCE [45] by  
 287 their official implementations. We leave the detailed settings in Appendix B.4.

### 288 4.3 Results

289 In Table 1, we compare the empirical results of NCI and corresponding baselines. On the NQ320k  
 290 dataset, the proposed NCI model outperforms all baselines by a large margin across four different  
 291 metrics. Compared with the state-of-the-art model, NCI improves 51.9% on Recall@1, 19.2% on  
 292 Recall@10, 6.2% on Recall@100, and 41.9% on MRR@100 relatively. It is worth noting that we  
 293 are the first to verify the superiority of deep text retrieval over traditional sparse and dense retrieval  
 294 methods. Previous deep text retrieval methods (*i.e.*, DSI and SEAL) obtain relatively poor results  
 295 even with a very large model size (e.g., T5-XXL). Consistent with previous studies, BM25 is an  
 296 efficient and effective baseline. It even outperforms BERT-based dense retrieval solutions, perhaps  
 297 owing to its capability to retrieve precise documents based on exact match. Further, we notice that  
 298 query generation plays a key role in boosting the retrieval performance. With query generation, the  
 299 BM25 + DocT5Query method achieves much higher performance than its vanilla version. ANCE  
 300 also achieves competitive performance after fine-tuned by the training pairs, but the performance is  
 301 far lower than our proposed NCI model. Moreover, the Recall@1 and MRR@100 metrics of NCI are  
 302 outstanding, indicating that more than 90% of the queries can be fulfilled without re-ranking on the  
 303 retrieved document list. This shows the potential of NCI to be served as an end-to-end solution that  
 304 replaces the entire index-retrieve-rank pipeline in traditional web search engines.

305 Furthermore, to study the effect of each component, we report ablation results on NQ320k dataset in  
 306 Table 2. In general, all five components are able to improve the performance of document retrieval,  
 307 which are detailed below.

308 **w/o query generation.** This configuration removes the query generation module for data augmenta-  
 309 tion. Remarkably, the query generation boosts the performance greatly. The result is aligned with  
 310 our expectation because training with the generated queries allows the model to be agnostic to the  
 311 semantic meaning of each documents. Besides, although training on  $\langle doc\text{-}content, docid \rangle$  pairs  
 312 like DSI [40] also make the model aware of the semantic meaning of each documents, we argue that

**Table 2:** Ablation Study on NQ320k retrieval task

Method	Recall@1	Recall@10	Recall@100	MRR@100
<b>Neural Corpus Indexer (Ours)</b>	<b>88.72</b>	<b>95.84</b>	<b>97.43</b>	<b>91.59</b>
w/o query generation	53.63	67.84	78.43	59.16
w/o PAWA decoder	87.01	95.27	97.18	90.79
w/o semantic id	87.22	95.34	97.25	90.85
w/o regularization	87.34	95.42	97.27	90.89
w/o constrained beam search	87.41	95.71	97.32	90.84

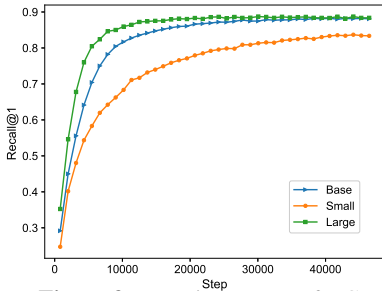
313 training with generated queries is able to avoid the distribution shift problem, which also benefit the  
 314 generalization performance.

315 **w/o PAWA decoder.** This configuration removes the adaptive decoder layer in Equation (4) and  
 316 leverages shared weights with token embedding for the linear classification layer. We notice that the  
 317 prefix-aware weight-adaptive decoder has a noticeable influence on the performance, which indicates  
 318 that, instead of borrowing the vanilla transformer decoder, it is necessary to design a tailored decoder  
 319 architecture for the task of semantic identifier generation.

320 **w/o semantic id.** This configuration replaces the semantic identifier of each document to a random  
 321 generated one. We find a relative drop in the model performance on all four metrics, demonstrating that  
 322 the semantic identifiers derived by hierarchical  $k$ -means have injected useful priors. We conjecture  
 323 that the performance enhancement would be more significant on a larger document corpus.

324 **w/o regularization.** There is a performance drop on all four metrics without using consistency-based  
 325 regularization loss. The reason is that the decoder network is prone to over-fitting. By making the  
 326 prediction results for two augmented versions of the decoder to be consistent, the decoder model  
 327 becomes more generalizable and resistant to over-fitting.

328 **w/o constrained beam search.** This configuration disables the validating constraint in beam search.  
 329 In other words, the decoder network does not have a tree-based prior structure. Instead, all tokens  
 330 in the vocabulary can be generated in each decoding step. We observe a performance drop on four  
 331 evaluation metrics. This indicates that it is difficult to remember all information of valid identifiers in  
 332 the network, and an explicit prior could be helpful for improving the quality of beam search.



333 **Figure 3:** Learning curves of NCI with different model capacities

Setting	Recall@1	Recall@10	Recall@100	MRR@100
#layer = 0	87.01	95.27	97.18	90.79
#layer = 1	88.54	95.62	97.16	91.44
#layer = 2	88.56	95.67	97.28	91.48
#layer = 4	88.65	95.72	<b>97.54</b>	91.51
#layer = 6	<b>88.72</b>	<b>95.84</b>	97.43	<b>91.59</b>
#layer = 8	85.31	94.17	96.34	89.25

**Table 3:** NCI with different number of layers in PAWA adapter and different regularization hyper-parameter  $\alpha$  in loss function

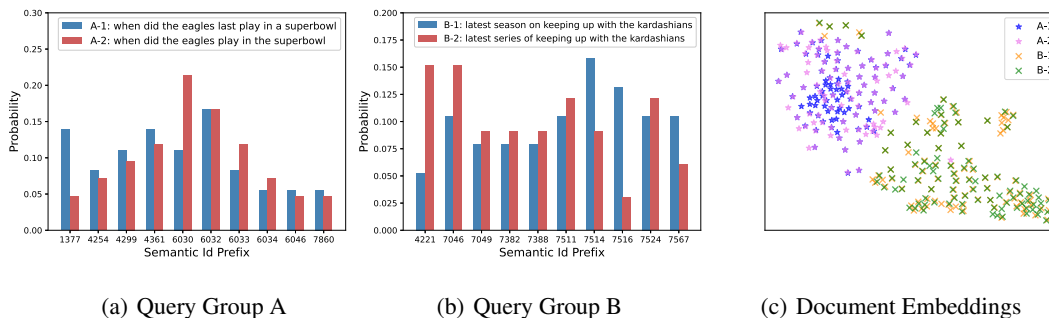
#### 334 4.4 Analysis

335 **Model capacity.** Figure 3 compares the learning curves of NCI with different model capacities,  
 336 which are identical to the small, base, and large settings of ordinary T5 [35]. We observe that with  
 337 the increase of model size, NCI convergences more quickly with fewer epochs. At convergence, the  
 338 small model achieves a relatively lower recall@1. Instead, both the base and large models achieve  
 339 similar results after sufficient training epochs. This implies that the model capacity has a critical  
 340 impact on the retrieval performance, and the capacity of base model seems to be enough to memorize  
 341 all documents in NQ320k dataset. For a larger corpus, one may need to increase the model size to  
 342 obtain satisfactory performance.

343 **Layer number of PAWA adapter.** We study the influence of the number of transformer layers in  
 344 the PAWA adapter and choose the layer number from  $\{0, 1, 2, 4, 6, 8\}$ . The results are summarized in  
 345 Table 3. We notice that with the increasing of layer number, *i.e.* from 0 to 6, the overall performance  
 346 is consistently improved under four metrics, except the Recall@100. But when the number of layer  
 347 achieves 8, the performance is dropped significantly. We attribute that to the overfitting caused by a  
 348 large PAWA adapter. Therefore, we adopt the design with 6 layer adapter in NCI.

349 **Retrieved documents and their semantics identifiers.** To verify the effectiveness of retrieval as  
 350 well as the semantic identifiers learned by hierarchical  $k$ -means, we analyze the retrieval results of  
 351 NCI for some exemplar queries. To illustrate, we select four queries denoted by  $A-1$ ,  $A-2$ ,  $B-1$  and  
 352  $B-2$ , where two queries inside the same group are semantically similar, and the queries in different  
 353 groups correspond to distinct topics. In Figure 4(a) and 4(b), we show the probabilities of retrieved  
 354 documents for each query in group  $A$  and  $B$  respectively. The digits along x-axis denote the four-bit  
 355 prefixes for semantic identifiers of retrieved documents, and the y-axis stands for their probabilities.





**Figure 4:** Analyses of retrieved documents with semantic identifiers. (a) The probabilities of retrieved documents for Query Group A; (b) Query Group B. (c) The t-SNE visualization of BERT-based document embeddings.

356 We notice that similar queries result in close document distributions, while dissimilar queries in  
 357 different groups result in un-overlapped document collections. In addition, the documents retrieved  
 358 by the same group of queries have close prefixes for the identifiers, *e.g.*, 6030, 6032, 6033, 6034 in  
 359 group A and 7511, 7514, 7516 in group B. Also, we visualize BERT-based document embeddings by  
 360 t-SNE [42] in Figure 4(c), in which each color represents the corresponding documents for a specific  
 361 query. As shown in the figure, these documents naturally form two clusters with respect to different  
 362 query groups. Thus, we conclude that the semantic document identifiers generated by the hierarchical  
 363  $k$ -means algorithm have positive effects on the retrieval performance.

364 **Efficiency Analysis.** We use an NVIDIA V100-32G GPU  
 365 to analyze the efficiency of NCI. As the inference speed is  
 366 influenced by both model capacity and beam size, we report  
 367 the latency and throughput measures for multiple settings  
 368 in Table 4. As NCI is an end-to-end retrieval method and  
 369 achieves competitive performance without re-ranking, the  
 370 latency and throughput are already affordable for some  
 371 near-real-time applications. **The latency of NCI is on par  
 372 with DSI and SEAL using the same model size and beam  
 373 size as all of them conduct beam search based on transformer decoders. BM25 is very efficient  
 374 (<100ms per query on CPU using an open-source implementation [2], but the recall metrics are much  
 375 lower.** Furthermore, we can leverage other techniques to improve the efficiency of NCI, which will  
 376 be discussed in the later section.

**Table 4: Efficiency analysis**

Model size	Beam size	Latency (ms)	Throughput (queries / s)
Small	10	78.46	58.48
Base	10	115.17	52.55
Large	10	188.60	43.39
Small	100	216.01	6.12
Base	100	269.31	5.62
Large	100	356.07	4.75

## 377 5 Limitation & Future Works

378 Despite the significant breakthrough, the current implementation of NCI still suffers from several  
 379 limitations before deployment in an industrial web search system. Firstly, it requires a much larger  
 380 model capacity for extending NCI to the web scale. Secondly, its inference speed needs to be  
 381 improved to serve online queries in real time. Thirdly, it is difficult to update the model-based index  
 382 when new documents are added to the system. In future works, we may tackle these problems from  
 383 four aspects. (1) The architecture of sparsely-gated Mixture of Expert (MoE) [38] can be employed  
 384 to enhance the model capacity. (2) Documents can be grouped into semantic clusters, and then NCI  
 385 is used to retrieve relevant cluster identifiers. In this way, all documents in relevant clusters can be  
 386 retrieved efficiently. (3) Model compression techniques, like weight quantization [22] and knowledge  
 387 distillation [20], can be further taken to speed up inference. (4) We plan to explore a hybrid solution  
 388 by building another index that serves new documents through traditional indexing algorithms.

## 389 6 Conclusion

390 In this work, we introduce a novel learning paradigm that unifies the learning and indexing stages by  
 391 an end-to-end deep neural network framework. The proposed Neural Corpus Indexer (NCI) retrieves  
 392 the identifiers of relevant documents directly for an input query, which can be optimized end-to-end  
 393 with augmented query-document pairs. To optimize the recall and ranking performance, we invent  
 394 the tailored prefix-aware weight-adaptive decoder. Empirically, we evaluate NCI on NQ320k dataset  
 395 and demonstrate its outstanding recall and MRR performance over state-of-the-art solutions.

## References

- 396  
397 [1] <https://github.com/castorini/docTTTTQuery>.
- 398 [2] <https://github.com/castorini/anserini>.
- 399 [3] Jon Louis Bentley. Multidimensional binary search trees used for associative searching. *Communications of the ACM*, 18(9):509–517, 1975.
- 400
- 401 [4] Michele Bevilacqua, Giuseppe Ottaviano, Patrick Lewis, Wen-tau Yih, Sebastian Riedel, and  
402 Fabio Petroni. Autoregressive search engines: Generating substrings as document identifiers.  
403 *arXiv preprint arXiv:2204.10628*, 2022.
- 404 [5] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
405 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
406 few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 407 [6] Wei-Cheng Chang, X Yu Felix, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. Pre-training  
408 tasks for embedding-based large-scale retrieval. In *International Conference on Learning*  
409 *Representations*, 2019.
- 410 [7] Qi Chen, Bing Zhao, Haidong Wang, Mingqin Li, Chuanjie Liu, Zengzhong Li, Mao Yang, and  
411 Jingdong Wang. Spann: Highly-efficient billion-scale approximate nearest neighbor search.  
412 *arXiv preprint arXiv:2111.08566*, 2021.
- 413 [8] Zhuyun Dai and Jamie Callan. Context-aware sentence/passage term importance estimation for  
414 first stage retrieval. *arXiv preprint arXiv:1910.10687*, 2019.
- 415 [9] Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S Mirrokni. Locality-sensitive hashing  
416 scheme based on p-stable distributions. In *Proceedings of the twentieth annual symposium on*  
417 *Computational geometry*, pages 253–262, 2004.
- 418 [10] Nicola De Cao, Wilker Aziz, and Ivan Titov. Highly parallel autoregressive entity linking with  
419 discriminative correction. *arXiv preprint arXiv:2109.03792*, 2021.
- 420 [11] Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. Autoregressive entity  
421 retrieval. *arXiv preprint arXiv:2010.00904*, 2020.
- 422 [12] Nicola De Cao, Ledell Wu, Kashyap Papat, Mikel Artetxe, Naman Goyal, Mikhail Plekhanov,  
423 Luke Zettlemoyer, Nicola Cancedda, Sebastian Riedel, and Fabio Petroni. Multilingual autore-  
424 gressive entity linking. *arXiv preprint arXiv:2103.12528*, 2021.
- 425 [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of  
426 deep bidirectional transformers for language understanding. In *North American Chapter of*  
427 *the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*,  
428 2019.
- 429 [14] Luyu Gao and Jamie Callan. Condenser: a pre-training architecture for dense retrieval. In  
430 *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*,  
431 pages 981–993, 2021.
- 432 [15] Luyu Gao, Zhuyun Dai, and Jamie Callan. Coil: Revisit exact lexical match in information  
433 retrieval with contextualized inverted list. In *Proceedings of the 2021 Conference of the*  
434 *North American Chapter of the Association for Computational Linguistics: Human Language*  
435 *Technologies*, pages 3030–3042, 2021.
- 436 [16] Weihao Gao, Xiangjun Fan, Chong Wang, Jiankai Sun, Kai Jia, Wenzhi Xiao, Ruofan Ding,  
437 Xingyan Bin, Hui Yang, and Xiaobing Liu. Deep retrieval: Learning a retrievable structure for  
438 large-scale recommendations. *arXiv preprint arXiv:2007.07203*, 2020.
- 439 [17] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. A deep relevance matching model  
440 for ad-hoc retrieval. In *Proceedings of the 25th ACM international on conference on information*  
441 *and knowledge management*, pages 55–64, 2016.

- 442 [18] Tonglei Guo, Jiafeng Guo, Yixing Fan, Yanyan Lan, Jun Xu, and Xueqi Cheng. A comparison  
443 between term-based and embedding-based methods for initial retrieval. In *China Conference on*  
444 *Information Retrieval*, pages 28–40. Springer, 2018.
- 445 [19] John A Hartigan and Manchek A Wong. Algorithm as 136: A k-means clustering algorithm.  
446 *Journal of the royal statistical society. series c (applied statistics)*, 28(1):100–108, 1979.
- 447 [20] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network.  
448 *arXiv preprint arXiv:1503.02531*, 2(7), 2015.
- 449 [21] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. Learning  
450 deep structured semantic models for web search using clickthrough data. In *Proceedings of*  
451 *the 22nd ACM international conference on Information & Knowledge Management*, pages  
452 2333–2338, 2013.
- 453 [22] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard,  
454 Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for  
455 efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer*  
456 *vision and pattern recognition*, pages 2704–2713, 2018.
- 457 [23] Suhas Jayaram Subramanya, Fnu Devvrit, Harsha Vardhan Simhadri, Ravishankar Krishnawamy,  
458 and Rohan Kadekodi. Diskann: Fast accurate billion-point nearest neighbor search on a single  
459 node. *Advances in Neural Information Processing Systems*, 32, 2019.
- 460 [24] Omar Khattab and Matei Zaharia. Colbert: Efficient and effective passage search via contextual-  
461 ized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR conference*  
462 *on research and development in Information Retrieval*, pages 39–48, 2020.
- 463 [25] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris  
464 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a  
465 benchmark for question answering research. *Transactions of the Association for Computational*  
466 *Linguistics*, 7:453–466, 2019.
- 467 [26] John Lafferty and Chengxiang Zhai. Document language models, query models, and risk  
468 minimization for information retrieval. In *Proceedings of the 24th annual international ACM*  
469 *SIGIR conference on Research and development in information retrieval*, pages 111–119, 2001.
- 470 [27] Canjia Li, Andrew Yates, Sean MacAvaney, Ben He, and Yingfei Sun. Parade: Passage  
471 representation aggregation for document reranking. *arXiv preprint arXiv:2008.09093*, 2020.
- 472 [28] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,  
473 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT  
474 pretraining approach. *CoRR*, abs/1907.11692, 2019.
- 475 [29] Wenhao Lu, Jian Jiao, and Ruofei Zhang. Twinbert: Distilling knowledge to twin-structured  
476 compressed bert models for large-scale retrieval. In *Proceedings of the 29th ACM International*  
477 *Conference on Information & Knowledge Management*, pages 2645–2652, 2020.
- 478 [30] Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. Sparse, dense, and  
479 attentional representations for text retrieval. *Transactions of the Association for Computational*  
480 *Linguistics*, 9:329–345, 2021.
- 481 [31] Yu A Malkov and Dmitry A Yashunin. Efficient and robust approximate nearest neighbor search  
482 using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and*  
483 *machine intelligence*, 42(4):824–836, 2018.
- 484 [32] Rodrigo Nogueira, Jimmy Lin, and AI Epistemic. From doc2query to docttttquery. *Online*  
485 *preprint*, 2019.
- 486 [33] Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. Document expansion by query  
487 prediction. *arXiv preprint arXiv:1904.08375*, 2019.

- 488 [34] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,  
489 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-  
490 learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830,  
491 2011.
- 492 [35] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
493 Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified  
494 text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2019.
- 495 [36] Stephen Robertson and Hugo Zaragoza. *The probabilistic relevance framework: BM25 and*  
496 *beyond*. Now Publishers Inc, 2009.
- 497 [37] Stephen E Robertson and Steve Walker. On relevance weights with little relevance information.  
498 In *Proceedings of the 20th annual international ACM SIGIR conference on Research and*  
499 *development in information retrieval*, pages 16–24, 1997.
- 500 [38] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,  
501 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts  
502 layer. *International Conference on Learning Representations (ICLR)*, 2017.
- 503 [39] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. Learning semantic  
504 representations using convolutional neural networks for web search. In *Proceedings of the 23rd*  
505 *international conference on world wide web*, pages 373–374, 2014.
- 506 [40] Yi Tay, Dara Bahri, Donald Metzler, Da-Cheng Juan, Zhe Zhao, and Che Zheng. Synthesizer:  
507 Rethinking self-attention in transformer models. *arXiv preprint arXiv:2005.00743*, 2020.
- 508 [41] Yi Tay, Vinh Q Tran, Mostafa Dehghani, Jianmo Ni, Dara Bahri, Harsh Mehta, Zhen Qin, Kai  
509 Hui, Zhe Zhao, Jai Gupta, et al. Transformer memory as a differentiable search index. *arXiv*  
510 *preprint arXiv:2202.06991*, 2022.
- 511 [42] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*  
512 *learning research*, 9(11), 2008.
- 513 [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
514 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Informa-*  
515 *tion Processing Systems (NeurIPS)*, pages 5998–6008, 2017.
- 516 [44] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony  
517 Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers:  
518 State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- 519 [45] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N Bennett, Junaid Ahmed,  
520 and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense  
521 text retrieval. In *International Conference on Learning Representations*, 2020.
- 522 [46] Wei Yang, Haotian Zhang, and Jimmy Lin. Simple applications of bert for ad hoc document  
523 retrieval. *arXiv preprint arXiv:1903.10972*, 2019.
- 524 [47] Hamed Zamani, Mostafa Dehghani, W Bruce Croft, Erik Learned-Miller, and Jaap Kamps. From  
525 neural re-ranking to neural ranking: Learning a sparse representation for inverted indexing.  
526 In *Proceedings of the 27th ACM international conference on information and knowledge*  
527 *management*, pages 497–506, 2018.
- 528 [48] Shengyao Zhuang, Hang Li, and G. Zuccon. Deep query likelihood model for information  
529 retrieval. In *ECIR*, 2021.

## 530 Checklist

531 The checklist follows the references. Please read the checklist guidelines carefully for information on  
532 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
533 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
534 the appropriate section of your paper or providing a brief inline description. For example:

- 535 • Did you include the license to the code and datasets? **[Yes]** See Section.
- 536 • Did you include the license to the code and datasets? **[No]** The code and the data are  
537 proprietary.
- 538 • Did you include the license to the code and datasets? **[N/A]**

539 Please do not modify the questions and only use the provided macros for your answers. Note that the  
540 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
541 block and only keep the Checklist section heading above along with the questions/answers below.

- 542 1. For all authors...
  - 543 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
544 contributions and scope? **[Yes]**
  - 545 (b) Did you describe the limitations of your work? **[Yes]** See Section 5.
  - 546 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See  
547 appendix.
  - 548 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
549 them? **[Yes]**
- 550 2. If you are including theoretical results...
  - 551 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
  - 552 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 553 3. If you ran experiments...
  - 554 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
555 mental results (either in the supplemental material or as a URL)? **[Yes]**
  - 556 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
557 were chosen)? **[Yes]** See Section 4.
  - 558 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
559 ments multiple times)? **[No]**
  - 560 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
561 of GPUs, internal cluster, or cloud provider)? **[Yes]** See Section 4.
- 562 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - 563 (a) If your work uses existing assets, did you cite the creators? **[Yes]** See Section 4.
  - 564 (b) Did you mention the license of the assets? **[Yes]** See appendix.
  - 565 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
  - 566 (d) Did you discuss whether and how consent was obtained from people whose data you're  
567 using/curating? **[Yes]**
  - 568 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
569 information or offensive content? **[Yes]**
- 570 5. If you used crowdsourcing or conducted research with human subjects...
  - 571 (a) Did you include the full text of instructions given to participants and screenshots, if  
572 applicable? **[N/A]**
  - 573 (b) Did you describe any potential participant risks, with links to Institutional Review  
574 Board (IRB) approvals, if applicable? **[N/A]**
  - 575 (c) Did you include the estimated hourly wage paid to participants and the total amount  
576 spent on participant compensation? **[N/A]**