

RIGHT: Retrieval-augmented Generation for Mainstream Hashtag Recommendation

Run-Ze Fan^{1,2}[0000-0002-8505-7756], Yixing Fan^{1,2}[0000-0003-4317-2702],
Jiangui Chen^{1,2}[0000-0002-6235-6526], Jiafeng Guo^{1,2*}[0000-0002-9509-8674],
Ruqing Zhang^{1,2}[0000-0003-4294-2541], and Xueqi Cheng^{1,2}[0000-0002-5201-8195]

¹ CAS Key Lab of Network Data Science and Technology, ICT, CAS

² University of Chinese Academy of Sciences, Beijing, China

fanrunze21s@ict.ac.cn

{fanyixing, chenjiangui18z, guojiafeng, zhangruqing, cxq}@ict.ac.cn

Abstract. Automatic mainstream hashtag recommendation aims to accurately provide users with concise and popular topical hashtags before publication. Generally, mainstream hashtag recommendation faces challenges in the comprehensive difficulty of newly posted tweets in response to new topics, and the accurate identification of mainstream hashtags beyond semantic correctness. However, previous retrieval-based methods based on a fixed predefined mainstream hashtag list excel in producing mainstream hashtags, but fail to understand the constant flow of up-to-date information. Conversely, generation-based methods demonstrate a superior ability to comprehend newly posted tweets, but their capacity is constrained to identifying mainstream hashtags without additional features. Inspired by the recent success of the retrieval-augmented technique, in this work, we attempt to adopt this framework to combine the advantages of both approaches. Meantime, with the help of the generator component, we could rethink how to further improve the quality of the retriever component at a low cost. Therefore, we propose *Retrieval-augmented Generative Mainstream HashTag Recommender (RIGHT)*, which consists of three components: (i) a retriever seeks relevant hashtags from the entire tweet-hashtags set; (ii) a selector enhances mainstream identification by introducing global signals; and (iii) a generator incorporates input tweets and selected hashtags to directly generate the desired hashtags. The experimental results show that our method achieves significant improvements over state-of-the-art baselines. Moreover, RIGHT can be easily integrated into large language models, improving the performance of ChatGPT by more than 10%. Code will be released at: <https://github.com/ict-bigdatalab/RIGHT>.

Keywords: Hashtag recommendation · Retrieval-augmented generation · Social media.

* Corresponding author

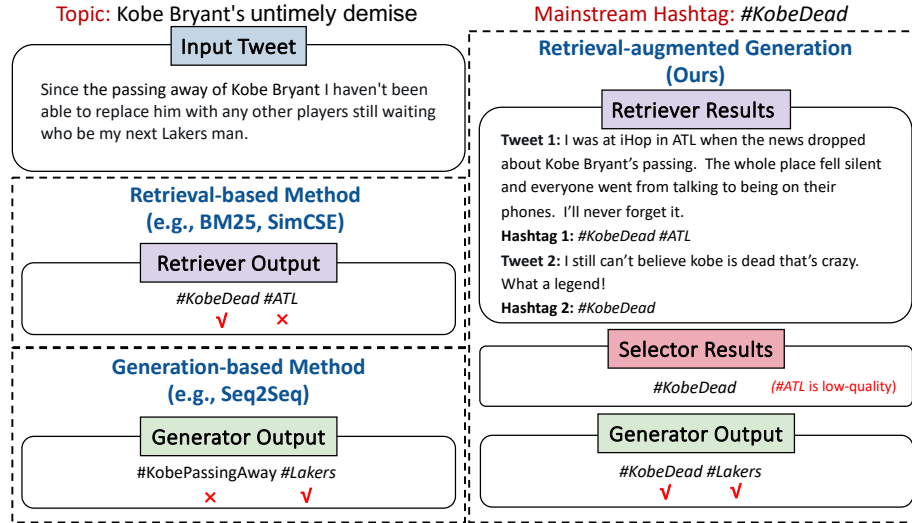


Fig. 1. Illustration of evaluating hashtag recommendation with different methods.

1 Introduction

Millions of user-generated microblogs flood Twitter daily, surpassing users' comprehension. To facilitate rapid and easy understanding, hashtags (e.g., *#ChatGPT*) are extensively used to convey central ideas and topics, which also enhance content visibility to reach a broader audience [24]. Such hashtags are commonly referred to as mainstream hashtags, denoting their status as not only the most prevalent hashtags but also semantically accurate. For instance, in the context of Kobe Bryant's untimely demise, both *#KobeDead* and *#KobeDeath* can be utilized, with both possessing accurate semantic meanings. However, the former holds more widespread usage and has achieved the distinction of being a mainstream hashtag.

To provide mainstream hashtags, two main challenges need to be addressed. First, comprehending a new tweet presents challenges primarily attributable to the absence of real-time information [42,50]. This is a direct consequence of the continuous emergence of numerous new tweets in response to new topics and events. Second, accurately identifying mainstream hashtags beyond semantic correctness remains a challenging task. The reason is that numerous hashtags could be used to describe a topic, but only a few are mainstream.

To address the above challenges, a considerable amount of work has been proposed, which could be divided into two research lines [10,13,43,15,42]. Retrieval-based methods retrieve hashtags from a fixed predefined mainstream hashtag list [43,15], which could alleviate the second problem. However, their ability to fully grasp the meaning of a newly posted tweet in response to emerging topics and events is constrained. Moreover, it is a considerable cost to maintain the predefined list [42]. In contrast, generation-based methods [42,50,30] demonstrate remarkable proficiency in comprehending new tweets and generating semantically accurate hashtags, owing to their substantial pretraining knowledge.

Nevertheless, they may encounter difficulties when it comes to identifying mainstream hashtags without enough mainstream information. As a result, the tweet might fail to be indexed by a mainstream hashtag on microblog services due to the tags’ unpopularity, weakening the recall rate of microblog searches. Inspired by the recent success of retrieval-augmented generation technique [1,37,21,12,48], therefore, we try to adapt this method to mainstream hashtags recommendation, utilizing the advantages of both retrieval and generation approaches.

Typically, retrieval-augmented techniques incorporate the results of the retriever, whether explicitly or implicitly, into the generator to enhance the quality of generation. Utilizing this framework, the introduction of a generator endowed with strong comprehensive capabilities might mitigate the dependency on the quality of the retriever [31]. Thus, we could rethink the trade-off between the quality and the cost of the retriever. Traditional retrieval-based methods rely on a predefined list of mainstream hashtags, which can ensure the quality of the retrieved information, but maintains such a list at a significant cost. To reduce the maintenance burden, we transform the small predefined list into a larger aggregation of existing tweet-hashtags pairs, which can be automatically collected and updated without manual cost. However, this approach carries the risk of introducing numerous low-quality hashtags due to the informal characteristics of social media content. Such hashtags have the potential to mislead the generator. As illustrated in Figure 1, both *#KobeDead* and *#ATL* are results of the retriever. Nonetheless, it is noteworthy that *#ATL* is low-quality, even though tweet 1 exhibits the highest degree of similarity with the input tweet. Consequently, it becomes imperative to further improve the quality of retrieved information without increasing the cost.

Therefore, in this study, we propose a *Retrieval-augmented Generative Mainstream HashTag Recommender (RIGHT)*, which combines the retriever and the generator by the retrieval-augmented technique with inserting a selector. Specifically, our method involves three components: 1) **Retriever** is utilized to acquire relevant hashtags. We retrieve the tweets most similar to the input from the tweet-hashtags corpus and obtain the corresponding hashtags set. 2) **Selector** is used to improve the capability of identifying mainstream hashtags. we incorporate three features, the similarity between the input tweet and the retrieved tweet and its hashtags, as well as the frequency of the hashtags, to enhance the mainstream information. 3) **Generator** is leveraged to provide strong semantic comprehension and the ability of hashtag generation. We concatenate the selected hashtags with the input tweet and feed it into the generator to obtain the desired hashtags. In this way, we can utilize not only the retriever and the selector to seek the mainstream hashtags but also the generator to produce the desired hashtags flexibly.

We conduct experiments on two large-scale datasets (i.e., English Twitter (THG) and Chinese Weibo (WHG)). Experimental results show that our method achieves significant improvements over state-of-the-art baselines. Moreover, as it can be easily incorporated in black-box language models, we also apply our framework to ChatGPT by zero-shot instruction learning, bringing a 12.7%

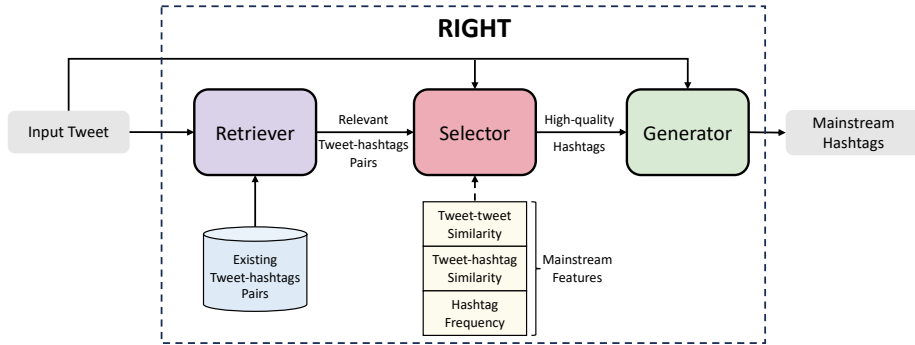


Fig. 2. Our RIGHT framework consists of a retriever, a selector, and a generator.

boost for THG and 18.3% for WHG in F1@1. Finally, to deeply understand this method, we present a detailed analysis.

2 Methodology

We propose *RetrIeval-augmented Generative Mainstream HashTag Recommender* (**RIGHT**), a simple yet effective framework for mainstream hashtag recommendation, which includes three components: retriever, selector, and generator.

Overall, we first utilize the retriever to retrieve the tweets most similar to the input from the existing tweet-hashtags corpus and obtain the corresponding hashtags set. Then, we adopt a selector to select the hashtags that are most probable mainstream from the retrieved labels using three signals. Finally, We concatenate the selected hashtags with the input tweet and feed it into the generator to obtain the desired hashtags. An overview of our method is shown in Figure 2.

2.1 Retriever

The goal of the retriever is to retrieve the top- N tweet-hashtags pairs on the same topic with the input tweet from the existing corpus, aiming to find relevant hashtags on the same topic.

Inspired by Wang et al. [41], we view the labeled training data as our corpus and index these as input-label pairs, i.e., $\mathcal{C} = \{(\tilde{t}_i, \tilde{H}_i)\}$. Then, given the input tweet t , the retrieval model \mathcal{R} matches it with all tweets in the corpus and returns the top- N most similar tweet-hashtags pairs together with their scores:

$$\{(\tilde{t}_1, \tilde{H}_1, \tilde{s}_1), \dots, (\tilde{t}_N, \tilde{H}_N, \tilde{s}_N)\} = \mathcal{R}(t|\mathcal{C}),$$

where we denote \tilde{s}_i as the similarity between t and the i -th retrieved tweet \tilde{t}_i . Each \tilde{H}_i consists of hashtags $\{\tilde{h}_1^i, \dots, \tilde{h}_{|\tilde{H}_i|}^i\}$. We report the results with sparse retrieval (e.g., BM25 [38]) and dense retrieval (e.g., SimCSE [11]) in experiments.

2.2 Selector

The goal of the selector is to filter the low-quality and non-mainstream hashtags existing in the results of the retriever (see Figure 1). We consider three mainstream features: the similarity between the input tweet and the retrieved tweet and its hashtags, as well as the frequency of the hashtags. Thus, we train a selector to compute the similarity between the tweet and the hashtags and propose a simple algorithm for hashtag ranking.

Training. The training data consists of positive samples and hard negative samples. Each hashtag labeled in a tweet can be viewed as a positive sample t^+ . However, a significant challenge lies in constructing hard negative samples (t^-) to facilitate the efficient selection of mainstream hashtags on the same topic by the selector. Inspired by BERT [9], we propose to create a hard negative sample by disturbing the labeled hashtag without changing the semantic meaning. Specifically, we randomly select a word to: (i) replace with its synonym 70% of the time; (ii) delete 10% of the time; (iii) swap with the adjacent word 10% of the time; (iv) insert a synonym after it 10% of the time. Thus, we obtain a training dataset $\{(t_i, t_i^+, t_i^-) | i = 1, \dots, N\}$. Finally, we utilize contrastive learning to train our selector by minimizing the following loss:

$$\mathcal{L}_{\mathcal{S}} = -\log \frac{e^{\text{sim}(\mathbf{h}_{t_i}, \mathbf{h}_{t_i^+})/\tau}}{\sum_{j=1}^L \left(e^{\text{sim}(\mathbf{h}_{t_i}, \mathbf{h}_{t_j^+})/\tau} + e^{\text{sim}(\mathbf{h}_{t_i}, \mathbf{h}_{t_j^-})/\tau} \right)}$$

where sim presents similarity, \mathbf{h}_t indicates the representation of t , L is mini-batch size, and τ is a temperature hyperparameter.

Inference. In the inference stage, we propose a simple algorithm for hashtag ranking. Given an input tweet t and the result of the retriever $\{(\tilde{t}_1, \tilde{H}_1, \tilde{s}_1), \dots, (\tilde{t}_N, \tilde{H}_N, \tilde{s}_N)\}$, we put retrieved hashtags into a set $\{\tilde{h}_1, \dots, \tilde{h}_M\}$ where we denote M as the number of different hashtags, and record the number of occurrences $\{f_1, \dots, f_M\}$ and the corresponding score of each retrieved hashtag $\{\tilde{s}_{i,1}, \dots, \tilde{s}_{i,f_i}\}$. Then, we match the input tweet t with all hashtags using the selector \mathcal{S} to obtain the similarity score between the tweet and all hashtags $\{\tilde{s}_1, \dots, \tilde{s}_M\}$:

$$\tilde{s}_m = \mathcal{S}(t, \tilde{h}_m). \quad (1)$$

Finally, we average the tweet-to-tweet similarity for each hashtag and add the similarity between the tweet and the hashtag. Since hashtags that occur more frequently are more likely to be mainstream, we magnify the sum of the similarity score and ranking score with a downscaled frequency:

$$s_i = \left(\frac{1}{f_i} \sum_{j=1}^{f_i} \tilde{s}_{i,j} \right) + \tilde{s}_i \times (1 + ((f_i - 1)/10)),$$

and sort the hashtags by the final score from largest to smallest to select the top- k hashtags $\{\tilde{h}_1, \dots, \tilde{h}_k\}$.

Table 1. Data statistics for the English Twitter hashtag generation (THG) dataset and the Chinese Weibo hashtag generation (WHG) dataset. # T-H pairs denotes the number of tweet-hashtags pairs. #AvgHashtags denotes the average number of hashtags in each tweet-hashtags pair. AvgTweetLen denotes the average length (token level) of all input tweets. AvgHashtagLen denotes the average length (token level) of all hashtags.

Dataset	THG Dataset			WHG Dataset		
	Train	Validation	Test	Train	Validation	Test
# T-H pairs	201444	11325	11328	307401	2000	2000
# AvgHashtags	4.1	4.1	4.1	1.0	1.0	1.0
AvgTweetLen	39.7	39.6	39.6	87.1	86.8	87.8
AvgHashtagLen	3.1	3.0	3.0	6.6	6.5	6.5

2.3 Generator

The goal of the generator is to generate the desired hashtags, given an input tweet t and the selected hashtags $\{\tilde{h}_1, \dots, \tilde{h}_k\}$. We concatenate the input tweet with the retrieved hashtags and separate each hashtag by a special token:

$$I = \langle t, \text{SEP1}, \tilde{h}_1, \text{SEP1}, \tilde{h}_2, \dots, \text{SEP1}, \tilde{h}_k \rangle,$$

and feed it into the generator \mathcal{G} , which will output a concatenated sequence O that includes the hashtags, with each hashtag separated by another token:

$$O = \langle h_1, \text{SEP2}, h_2, \dots, \text{SEP2}, h_{|H|} \rangle.$$

By easily splitting by the special token, we would obtain the hashtag list $H = \{h_1, h_2, \dots, h_{|H|}\}$.

The generative model could be Transformer-based encoder-decoder architecture (e.g., T5 [35], BART [25]) or decoder-only architecture (e.g., a series of GPT [33,34,2]). Thus, the training stage focuses on the finetuning of generative models by minimizing the cross-entropy loss:

$$\mathcal{L}_{\mathcal{G}} = \sum_{(I,O) \in \mathcal{D}} -\log p(O|I; \theta_{\mathcal{G}}),$$

where I is the input sequence consisting of the input tweet and the selected hashtags, O is the output sequence consisting of the desired hashtags, and $\theta_{\mathcal{G}}$ is the parameters of the generator.

3 Experiments

3.1 Experimental Setup

Datasets. Our experiments are conducted on two large-scale datasets, which were crawled from official media and influencers of social media [30]. The details are shown in table 1.

- **THG:** The English Twitter hashtag generation (THG) dataset has been crawled from official Twitter sources, encompassing organizations, media outlets, and other authenticated users, with the primary objective of acquiring tweets of superior quality.
- **WHG:** The Chinese Weibo hashtag generation (WHG) dataset has been acquired through the systematic extraction of microblogs from Weibo, encompassing notable sources including *People’s Daily*, *People.cn*, *Economic Observe press*, *Xinlang Sports*, and various other accounts boasting over 5 million followers. These accounts span diverse domains, encompassing politics, economics, military affairs, sports, and more.

We use the training datasets as our retrieval corpus.

Evaluation Metric. Following previous work [42,30], we utilize ROUGE metrics and F1 scores at K as our evaluation metric. The average ROUGE score measures the overlap between the generated sequence of hashtags (excluding special tokens) and the reference sequence, including ROUGE-1, ROUGE-2, and ROUGE-L. For F1 scores at K , different K values result in a similar trend, so only F1@1 and F1@5 are reported. We report results on the test dataset. Noticeably, for the WHG dataset, where input posts have only one hashtag, F1@1 and F1@5 are identical, so we only report F1@1 for this dataset.

Implementation Details. Our implementation details of the retriever, selector, and generator are the following:

- **For Retriever**, we utilize BM25 [38] and SimCSE [11] (i.e., RoBERTa-Large [29] for THG and Bert-Base-Chinese [9] for WHG) as our retrievers. Following Gao et al. [11], we train our model for 3 epochs with a learning rate of $1e-5$. The hyperparameter of N is set to 10, and the batch size is 6 per device.
- **For Selector**, we use the training datasets from THG and WHG to construct our hard negative samples, which are subsequently employed for training our selectors in both English and Chinese independently. We utilize RoBERTa-Large for THG and Bert-Base-Chinese for WHG. The temperature τ is 0.05 and other hyperparameters are the same as the retriever.
- **For Generator**, we fine-tune a T5-base [35] for THG and a mT5-small [44] for WHG and use Adam [22] as an optimizer. We set the weight decay and batch size as $1e-5$ and 16 and grid-search the learning rate, training epochs, and the number of concatenated hashtags k from $\{3e-4, 1e-4, 5e-5\}$, $\{5, 10\}$, and $\{1, 3, 5, 7, 9\}$ respectively. The maximum length is 180 for T5-base and 256 for mT5-small. The special token SEP1 and SEP2 are `< extra_id_0 >` and `< extra_id_1 >` respectively.

All models are trained on four NVIDIA Tesla K80.

Table 2. The prompt used for ChatGPT. The Chinese version is its translation.

Baseline	Instruction
ChatGPT	I want you to act as a hashtag annotator. I will provide you a tweet and your role is to annotate the relevant hashtag. You should use the related knowledge and find the topic. I want you only reply the hashtags segmented by “#” and nothing else, do not write explanations. I want you segment the word in a hashtag by space. My first tweet is <code>{Input Tweet}</code> .
RIGHT^{ChatGPT}	I want you to act as a hashtag annotator. I will provide you with a tweet, and your role is to annotate the relevant hashtag. Using your related knowledge, you should identify the topic and reply with only the hashtags segmented by “#”, without any explanations. Make sure to capitalize the first letter of the word. Make sure to split every word in a hashtag by a space. There are some potential hashtags: <code>{{Retrieved Top-k Hashtags}}</code> . You can decide whether use the part of them or not. My first tweet is <code>{Input Tweet}</code> .

Baselines. Our baselines consist of retrieval-based methods, generation-based methods, and retrieval-augmented generative methods:

- **Retrieval-based methods:** Following Mao et al. [30], we construct the predefined hashtags list from all hashtags in the training datasets and select top-4 hashtags for THG and top-1 for WHG according to the average number of hashtags in each data item. We apply BM25 and SimCSE to the hashtag recommendation: (i) **BM25** [38] is a traditional strong sparse retrieval based on term matching. (ii) **SimCSE** [11] is a representative dense retriever, which applies a simple contrastive learning framework to present sentence embeddings on semantic textual similarity tasks. We fine-tuned SimCSE by constructing positive samples and hard negative samples from BM25.
- **Generation-based methods:** We consider three predominant generative methods: (i) **ChatGPT** is a powerful large language model to execute various NLP tasks [23]. Specifically, we adopt `gpt-3.5-turbo` and instruction zero-shot learning to evaluate our task (Prompts are shown in Table 2). (ii) **SEGTRM Soft** [30] is the previous SOTA on our datasets, an end-to-end generative method segments selection-based deep transformer. (iii) **Seq2Seq** [42] is the first generation-based method for hashtag recommendation. Due to the unreality to assume the existence of conversations before publishing the tweet [50] and the lack of conversation contexts, we reimplement the Seq2Seq model on the pretrained language model (T5-base for THG and mT5-small for WHG) to formulate this task to a seq-to-seq paradigm.
- **Retrieval-augmented Generative Methods (Ours):** We apply our retrieval-augmented framework to ChatGPT by incorporating the retrieval results into the instruction to prompt the model to generate mainstream hashtags, denoted as **RIGHT^{ChatGPT}** (Prompts are shown in Table 2). We

Table 3. Main results (%) on the THG and WHG datasets. Bold and underline indicate the best and second method respectively. We denote ROUGE as RG. * indicates statistically significant improvements over all baselines (p-value < 0.05).

Model	THG					WHG			
	RG-1	RG-2	RG-L	F1@1	F1@5	RG-1	RG-2	RG-L	F1@1
<i>Retrieval-based Methods</i>									
BM25	16.23	4.17	15.11	5.92	9.84	61.98	58.76	61.81	48.20
SimCSE	28.43	10.34	26.38	12.40	15.15	59.71	55.81	59.54	47.65
<i>Generation-based Methods</i>									
ChatGPT	44.60	27.67	39.29	9.72	26.08	32.27	24.54	31.80	7.9
SEGTRM Soft	51.18	37.15	47.05	27.17	29.02	55.51	51.28	54.30	30.72
Seq2Seq	59.90	41.39	59.15	29.75	41.71	66.64	61.71	66.39	48.60
<i>Retrieval-augmented Generative Methods (Ours)</i>									
RIGHT ^{ChatGPT}	47.54	25.63	44.47	22.39	31.09	48.17	41.51	47.75	26.15
RIGHT _{BM25}	<u>61.60</u>	<u>43.77</u>	<u>60.85</u>	<u>30.27</u>	<u>42.98</u>	70.62*	66.12*	70.35*	53.85*
RIGHT _{SimCSE}	62.11*	43.86*	61.39*	30.58*	43.23*	<u>68.84</u>	<u>64.19</u>	<u>68.56</u>	<u>51.50</u>

only use the best retriever on the datasets (i.e., SimCSE for THG and BM25 for WHG), due to the high cost of ChatGPT. Moreover, we use BM25 and SimCSE as our retriever of RIGHT, denoted them as **RIGHT**_{BM25} and **RIGHT**_{SimCSE}.

3.2 Main Results

As shown in Table 3, we can observe that:

1. Among the retrieval-based methods, the performance of SimCSE outperforms BM25 in THG, while BM25 demonstrates superior performance in WHG. This difference may be attributed to that English hashtags tend to be concise summaries, while Chinese hashtags often comprise small sentences extracted directly from the input text. Consequently, dense retrieval approaches utilizing semantic matching may be more suitable for English datasets, while Chinese datasets may benefit more from sparse retrieval techniques based on term matching.
2. Among the generation-based methods, Seq2Seq performs well on both datasets, potentially attributable to the utilization of mainstream hashtag knowledge from the training dataset during fine-tuning. However, ChatGPT lags behind other generation methods, suggesting a deficiency in mainstream hashtag knowledge despite its vast repository of general knowledge.
3. Among the retrieval-augmented generative methods, retrieval augmentation brings the performance of baselines to a new level, demonstrating the effectiveness of our method. For ChatGPT, retrieval augmentation boosts F1@1 performance by 12.67% for THG and 18.25% for WHG, indicating the substantial value of mainstream hashtag knowledge. For RIGHT, both sparse

Table 4. Ablation study results on the THG datasets. Bold indicates the best method.

Model	ROUGE-1	ROUGE-2	ROUGE-L	F1@1	F1@5
RIGHT	62.11	43.86	61.39	30.58	43.23
w/o Retriever	59.91	41.63	59.23	29.66	41.70
w/o Selector	60.49	42.06	59.76	30.22	41.95
w/o Generator	36.24	16.02	32.86	24.61	26.73

and dense retrievers show the potential to enhance performance compared with Seq2Seq. Specifically, SimCSE is particularly effective for THG, while BM25 performs better for WHG. The reason could be attributed to the superiority in the performance of retrieval-based methods is directly proportional to the enhancement of the retrieval augmentation. Moreover, different retrievers excel in different scenarios, emphasizing the importance of the careful selection of retrievers based on specific use cases. The performance of the retrieval-based approach serves as a preliminary guide for informed decision-making.

Overall, our method shows robustness across various scenarios, whether applied with a fine-tuned generation model or a large black box language model. Regardless of the retrieval approach used, our method consistently improves performance.

3.3 Analysis

Ablation Study. We conduct an ablation study to explore the impact of each component in RIGHT on THG: 1) **w/o Retriever:** We remove the retriever and randomly concatenate k hashtags from the training dataset with the input tweet. 2) **w/o Selector:** We remove the selector and directly use the top- k hashtags from the retriever’s results by the similarity between the input tweet and the retrieved tweet. 3) **w/o Generator:** We remove the generator and output the top-4 hashtags produced by the selector. Table 4 presents the results, indicating that:

1. The performance improvement is considerable through the integration of the retriever, confirming that the incorporation of mainstream hashtag knowledge indeed facilitates accurate hashtag selection.
2. Without the selector, the performance gains are limited, indicating that simply being on the same topic is insufficient. It is crucial to identify and incorporate mainstream hashtags.
3. The generator is crucial in RIGHT, emphasizing the significant impact of semantic comprehension on performance. In contrast to the retrieval-based approaches, it is more powerful to directly output the hashtags in the tweet-hashtags pair that are most similar to the input.

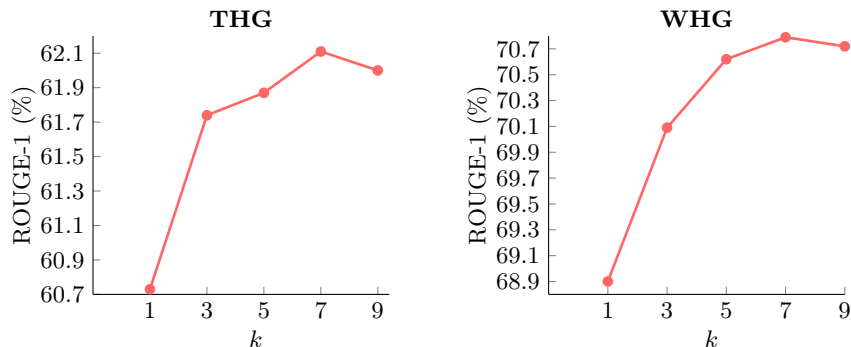


Fig. 3. Our Rouge-1 results in the different number of augmented hashtags (i.e., $k = 1, 3, 5, 7, 9$).

Impact of the Number of Augmented Hashtags. To explore the impact of the number of concatenated hashtags with the input tweet, we conduct a series of experiments. Specifically, we concatenate various top- k ($k = 1, 3, 5, 7, 9$) with the input tweet for the THG and the WHG datasets. Figure 3 demonstrates the Rouge-1 results (other metrics show the same trends), showing that:

1. Retrieval augmentation aids in improving the performance when a sufficient number is considered, suggesting that augmenting more hashtags increases the probability of covering mainstream hashtags and makes the generator more robust in the presence of mismatches from certain hashtags.
2. Upon reaching a certain threshold of the number of augmented hashtags (i.e., $k = 7$), the performance converges, suggesting that the majority of mainstream hashtags might have already been augmented.

Case Study. To validate the successful recall of mainstream hashtags and the potential for further improvement, we analyze the successful and unsuccessful cases and present a representative case study in Table 5. We conclude that:

1. Some retrieved tweets share the same topic as the input tweets but have sub-par labeled hashtags (e.g., “fx logix” in Table 5). Fortunately, our generator demonstrates the capability to disregard these irrelevant hashtags.
2. Some retrieved tweets are partially relevant to the input tweet. Although the retrieved hashtags align well with the topic of the retrieved tweet, it is not highly pertinent to the primary topic of the input tweet (e.g., “azure” in the retrieved hashtags). Nonetheless, our generator can filter out these irrelevant hashtags.
3. The generation model produces a semantically accurate but non-mainstream hashtag “windows virtual desktop” by directly copying the original word from the input tweet due to its limited knowledge of mainstream hashtags. However, our retriever and selector effectively identify the corresponding mainstream hashtag in its abbreviated form “wvd”. Our RIGHT successfully

Table 5. An example from the THG test set. Correct results are marked bold.

Input:	Geeks guide to microsoft teams optimization with windows virtual desktop citrix.
Label:	microsoft; windows; citrix; wvd
Retriever & Selector:	citrix; wvd ; fs logix; v mware; azure; aws; microsoft
Seq2Seq:	microsoft ; windows virtual desktop; citrix ; vdi
RIGHT:	citrix; wvd

replaced the original hashtag with the mainstream hashtag, indicating the effectiveness of retrieval augmentation.

- Seq2Seq generates certain hashtags that are also retrieved by the retriever, while RIGHT fails to generate them (e.g., “microsoft”). These cases constitute less than 1% of the total. We speculate that this discrepancy may be due to the selector placing the correct hashtags toward the end of the list, leading to reduced confidence and subsequent non-adoption by the generator.

4 Related Works

Our work mainly builds on two streams of previous work: hashtag recommendation and retrieval-augmented generation.

4.1 Mainstream hashtag recommendation

Mainstream hashtag recommendation aims to provide users with short topical and popular tags representing the main ideas of their tweets before publication. Three primary methods have been proposed for this task [24]:

- 1) **Keyphrase extraction method** formulates this task as keyphrases extraction from source posts [14,47,49], which fails to produce hashtags that do not appear in the microblog posts while large freedom is allowed for users to write whatever hashtags they like. The performance of this method is much lower than other methods.
- 2) **Retrieval-based method** aims to retrieve from a predefined hashtag list [43,15,19,46], which is limited to generating only the hashtags that are included in the list. In reality, a wide range of hashtags can be created every day, resulting impossibility to be covered by a fixed list and the difficulty to maintain the list.
- 3) **Generation-based method** was proposed [42,50,30,32] to overcome the aforementioned challenges, which formulates the task as a sequence-to-sequence generation paradigm, allowing for the creation of a wider range of hashtags that better capture the main ideas of the microblog post. However, previous studies pay limited attention to mainstream hashtags. Consequently, even though it produces semantically correct tags, the tweet might fail to be indexed by a mainstream hashtag on microblog services due to tags’ unpopularity, thus weakening the recall rate of microblog searches.

To the best of our knowledge, we are the first to alleviate this issue by combining retrieval and generation methods. Meanwhile, we improve the quality of the retriever at a low cost.

4.2 Retrieval-augmented Generation

The retrieval-augmented generation represents a novel paradigm that merges pre-trained generative models with information retrieval techniques [1,26]. Previous research in this field primarily has focused on introducing external knowledge to address knowledge-intensive tasks [4,45,18,27,39,40,17,8,6,7] and utilizing similar data to enhance the model performance across various natural language processing (NLP) tasks, including image captioning [37], keyphrase generation [21,12], named entity recognition [48,5], and others. Recently, this technique has also been used in large language models to alleviate issues like factual hallucination [3,36,20], knowledge out-dating [16], and the lack of domain-specific expertise [28].

Notably, we adopt this framework to the mainstream hashtag recommendation task, by introducing a selector combining global signals to improve mainstream identification.

5 Conclusion

In this study, we have proposed a simple yet effective retrieval-augmented generative recommender, designed to utilize the advantage of retrieval and generation methods for mainstream hashtag recommendation. To improve the quality of the retriever’s results at a low cost, we have integrated a selector module into the conventional retrieval-augmented framework. Specifically, the retriever’s role is to find relevant hashtags on the same topic, the selector is employed to enhance the identification of mainstream hashtags, and the generator is responsible for combining input tweets and selected hashtags to generate desired hashtags. We have conducted extension experiments using two extensive datasets to validate the effectiveness of our approach.

In future work, it is valuable to explore optimal strategies for combining the retrieval-based method with the generation-based method, as well as developing a co-training approach that jointly refines the three components.

Acknowledgements

This work was funded by the National Natural Science Foundation of China (NSFC) under Grants No. 62372431, and 62006218, the Youth Innovation Promotion Association CAS under Grants No. 2021100, the project under Grants No. 2023YFA1011602, JCKY2022130C039 and 2021QY1701, and the Lenovo-CAS Joint Lab Youth Scientist Project. All content represents the opinion of the authors, which is not necessarily shared or endorsed by their respective employers and/or sponsors.

References

1. Asai, A., Min, S., Zhong, Z., Chen, D.: Retrieval-based language models and applications. In: ACL. pp. 41–46. Toronto, Canada (Jul 2023). <https://doi.org/10.18653/v1/2023.acl-tutorials.6>, <https://aclanthology.org/2023.acl-tutorials.6>
2. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D.: Language models are few-shot learners. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) NeurIPS (2020), <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfbcb4967418bfb8ac142f64a-Abstract.html>
3. Cao, M., Dong, Y., Wu, J., Cheung, J.C.K.: Factual error correction for abstractive summarization models. In: EMNLP (Nov 2020). <https://doi.org/10.18653/v1/2020.emnlp-main.506>, <https://aclanthology.org/2020.emnlp-main.506>
4. Chen, D., Fisch, A., Weston, J., Bordes, A.: Reading Wikipedia to answer open-domain questions. In: ACL (Jul 2017). <https://doi.org/10.18653/v1/P17-1171>, <https://aclanthology.org/P17-1171>
5. Chen, J., Zhang, R., Guo, J., Fan, Y., Cheng, X.: GERE: generative evidence retrieval for fact verification. In: SIGIR 2022. pp. 2184–2189. ACM (2022). <https://doi.org/10.1145/3477495.3531827>, <https://doi.org/10.1145/3477495.3531827>
6. Chen, J., Zhang, R., Guo, J., Liu, Y., Fan, Y., Cheng, X.: Corpusbrain: Pre-train a generative retrieval model for knowledge-intensive language tasks. In: CIKM 2022. pp. 191–200. ACM (2022). <https://doi.org/10.1145/3511808.3557271>, <https://doi.org/10.1145/3511808.3557271>
7. Chen, J., Zhang, R., Guo, J., de Rijke, M., Chen, W., Fan, Y., Cheng, X.: Continual learning for generative retrieval over dynamic corpora. In: CIKM 2023. pp. 306–315. ACM (2023). <https://doi.org/10.1145/3583780.3614821>, <https://doi.org/10.1145/3583780.3614821>
8. Chen, J., Zhang, R., Guo, J., de Rijke, M., Liu, Y., Fan, Y., Cheng, X.: A unified generative retriever for knowledge-intensive language tasks via prompt learning. In: SIGIR 2023. pp. 1448–1457. ACM (2023). <https://doi.org/10.1145/3539618.3591631>, <https://doi.org/10.1145/3539618.3591631>
9. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: NAACL (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>, <https://aclanthology.org/N19-1423>
10. Ding, Z., Zhang, Q., Huang, X.: Automatic hashtag recommendation for microblogs using topic-specific translation model. In: COLING (Dec 2012), <https://aclanthology.org/C12-2027>
11. Gao, T., Yao, X., Chen, D.: SimCSE: Simple contrastive learning of sentence embeddings. In: EMNLP (Nov 2021). <https://doi.org/10.18653/v1/2021.emnlp-main.552>, <https://aclanthology.org/2021.emnlp-main.552>
12. Gao, Y., Yin, Q., Li, Z., Meng, R., Zhao, T., Yin, B., King, I., Lyu, M.: Retrieval-augmented multilingual keyphrase generation with retriever-generator iterative training. In: NAACL (2022). <https://doi.org/10.18653/v1/2022.findings-naacl.92>, <https://aclanthology.org/2022.findings-naacl.92>

13. Gong, Y., Zhang, Q., Huang, X.: Hashtag recommendation using Dirichlet process mixture models incorporating types of hashtags. In: EMNLP (2015). <https://doi.org/10.18653/v1/D15-1046>, <https://aclanthology.org/D15-1046>
14. Gong, Y., Zhang, Q., Huang, X.: Hashtag recommendation using dirichlet process mixture models incorporating types of hashtags. In: Màrquez, L., Callison-Burch, C., Su, J., Pighin, D., Marton, Y. (eds.) EMNLP (2015). <https://doi.org/10.18653/v1/d15-1046>, <https://doi.org/10.18653/v1/d15-1046>
15. Gong, Y., Zhang, Q.: Hashtag recommendation using attention-based convolutional neural network. In: Kambhampati, S. (ed.) IJCAI (2016), <http://www.ijcai.org/Abstract/16/395>
16. He, H., Zhang, H., Roth, D.: Rethinking with retrieval: Faithful large language model inference. arXiv preprint (2022), <https://arxiv.org/pdf/2301.00303.pdf>
17. He, S., Fan, R.Z., Ding, L., Shen, L., Zhou, T., Tao, D.: Mera: Merging pretrained adapters for few-shot learning. arXiv preprint arXiv:2308.15982 (2023), <https://arxiv.org/abs/2308.15982>
18. He, S., Fan, R.Z., Ding, L., Shen, L., Zhou, T., Tao, D.: Merging experts into one: Improving computational efficiency of mixture of experts. In: Bouamor, H., Pino, J., Bali, K. (eds.) Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing. pp. 14685–14691. Association for Computational Linguistics, Singapore (Dec 2023), <https://aclanthology.org/2023.emnlp-main.907>
19. Huang, H., Zhang, Q., Gong, Y., Huang, X.: Hashtag recommendation using end-to-end memory networks with hierarchical attention. In: COLING (Dec 2016), <https://aclanthology.org/C16-1090>
20. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y.J., Madotto, A., Fung, P.: Survey of hallucination in natural language generation. ACM Comput. Surv. (2023). <https://doi.org/10.1145/3571730>, <https://doi.org/10.1145/3571730>
21. Kim, J., Jeong, M., Choi, S., Hwang, S.w.: Structure-augmented keyphrase generation. In: EMNLP (2021). <https://doi.org/10.18653/v1/2021.emnlp-main.209>, <https://aclanthology.org/2021.emnlp-main.209>
22. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Bengio, Y., LeCun, Y. (eds.) ICLR (2015), <http://arxiv.org/abs/1412.6980>
23. Kocoń, J., Cichecki, I., Kaszyca, O., Kochanek, M., Szydło, D., Baran, J., Bielaniec, J., Gruza, M., Janz, A., Kanclerz, K., Kocoń, A., Koptyra, B., Mieszczenko-Kowszewicz, W., Miłkowski, P., Oleksy, M., Piasecki, M., Łukasz Radliński, Wojtasik, K., Woźniak, S., Kazienko, P.: Chatgpt: Jack of all trades, master of none. arXiv preprint (2023), <https://arxiv.org/pdf/2302.10724.pdf>
24. Kwak, H., Lee, C., Park, H., Moon, S.: What is twitter, a social network or a news media? In: WWW (2010), https://ink.library.smu.edu.sg/cgi/viewcontent.cgi?article=7104&context=sis_research
25. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L.: BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In: Jurafsky, D., Chai, J., Schluter, N., Tetreault, J.R. (eds.) ACL (2020). <https://doi.org/10.18653/v1/2020.acl-main.703>, <https://doi.org/10.18653/v1/2020.acl-main.703>
26. Lewis, P.S.H., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W., Rocktäschel, T., Riedel, S., Kiela, D.: Retrieval-augmented generation for knowledge-intensive NLP tasks.

- In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) NeurIPS (2020), <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>
27. Li, J., Sun, S., Yuan, W., Fan, R.Z., Zhao, H., Liu, P.: Generative judge for evaluating alignment. arXiv preprint arXiv:2310.05470 (2023), <https://arxiv.org/abs/2310.05470>
 28. Li, X., Zhu, X., Ma, Z., Liu, X., Shah, S.: Are chatgpt and gpt-4 general-purpose solvers for financial text analytics? an examination on several typical tasks. arXiv preprint (2023), <https://arxiv.org/pdf/2305.05862.pdf>
 29. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. arXiv preprint (2019), <https://arxiv.org/pdf/1907.11692.pdf>
 30. Mao, Q., Li, X., Liu, B., Guo, S., Hao, P., Li, J., Wang, L.: Attend and select: A segment selective transformer for microblog hashtag generation. Knowledge-Based Systems (2022). <https://doi.org/https://doi.org/10.1016/j.knosys.2022.109581>, <https://www.sciencedirect.com/science/article/pii/S0950705122007973>
 31. Mialon, G., Dessì, R., Lomeli, M., Nalmpantis, C., Pasunuru, R., Raileanu, R., Rozière, B., Schick, T., Dwivedi-Yu, J., Celikyilmaz, A., et al.: Augmented language models: a survey. arXiv preprint (2023), <https://arxiv.org/pdf/2302.07842.pdf>
 32. Ni, S., Bi, K., Guo, J., Cheng, X.: A comparative study of training objectives for clarification facet generation. In: Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region. pp. 1–10 (2023)
 33. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al.: Improving language understanding by generative pre-training. preprint (2018), https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf
 34. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al.: Language models are unsupervised multitask learners. preprint (2019), https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf
 35. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. J. Mach. Learn. Res. (2020), <http://jmlr.org/papers/v21/20-074.html>
 36. Raunak, V., Menezes, A., Junczys-Dowmunt, M.: The curious case of hallucinations in neural machine translation. In: ACL (Jun 2021). <https://doi.org/10.18653/v1/2021.naacl-main.92>, <https://aclanthology.org/2021.naacl-main.92>
 37. Rita Ramos, Desmond Elliott, B.M.: Retrieval-augmented image captioning. In: EACL (2023), <https://arxiv.org/abs/2302.08268>
 38. Robertson, S.E., Walker, S.: Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In: Croft, W.B., van Rijsbergen, C.J. (eds.) SIGIR (1994). https://doi.org/10.1007/978-1-4471-2099-5_24, https://doi.org/10.1007/978-1-4471-2099-5_24
 39. Wang, S., Gan, T., Liu, Y., Wu, J., Cheng, Y., Nie, L.: Micro-influencer recommendation by multi-perspective account representation learning. IEEE Transactions on Multimedia (2022)
 40. Wang, S., Gan, T., Liu, Y., Zhang, L., Wu, J., Nie, L.: Discover micro-influencers for brands via better understanding. IEEE Transactions on Multimedia **24**, 2595–2605 (2021)

41. Wang, S., Xu, Y., Fang, Y., Liu, Y., Sun, S., Xu, R., Zhu, C., Zeng, M.: Training data is more valuable than you think: A simple and effective method by retrieving from training data. In: ACL (May 2022). <https://doi.org/10.18653/v1/2022.acl-long.226>, <https://aclanthology.org/2022.acl-long.226>
42. Wang, Y., Li, J., King, I., Lyu, M.R., Shi, S.: Microblog hashtag generation via encoding conversation contexts. In: NAACL (Jun 2019). <https://doi.org/10.18653/v1/N19-1164>, <https://aclanthology.org/N19-1164>
43. Weston, J., Chopra, S., Adams, K.: #TagSpace: Semantic embeddings from hashtags. In: EMNLP (Oct 2014). <https://doi.org/10.3115/v1/D14-1194>, <https://aclanthology.org/D14-1194>
44. Xue, L., Constant, N., Roberts, A., Kale, M., Al-Rfou, R., Siddhant, A., Barua, A., Raffel, C.: mt5: A massively multilingual pre-trained text-to-text transformer. In: Toutanova, K., Rumshisky, A., Zettlemoyer, L., Hakkani-Tür, D., Beltagy, I., Bethard, S., Cotterell, R., Chakraborty, T., Zhou, Y. (eds.) NAACL (2021). <https://doi.org/10.18653/v1/2021.naacl-main.41>, <https://doi.org/10.18653/v1/2021.naacl-main.41>
45. Zhang, H., Zhang, R., Guo, J., de Rijke, M., Fan, Y., Cheng, X.: From relevance to utility: Evidence retrieval with feedback for fact verification. In: Bouamor, H., Pino, J., Bali, K. (eds.) Findings of the Association for Computational Linguistics: EMNLP 2023. pp. 6373–6384. Association for Computational Linguistics, Singapore (Dec 2023), <https://aclanthology.org/2023.findings-emnlp.422>
46. Zhang, Q., Wang, J., Huang, H., Huang, X., Gong, Y.: Hashtag recommendation for multimodal microblog using co-attention network. In: Sierra, C. (ed.) IJCAI (2017). <https://doi.org/10.24963/ijcai.2017/478>, <https://doi.org/10.24963/ijcai.2017/478>
47. Zhang, Q., Wang, Y., Gong, Y., Huang, X.: Keyphrase extraction using deep recurrent neural networks on Twitter. In: EMNLP (Nov 2016). <https://doi.org/10.18653/v1/D16-1080>, <https://aclanthology.org/D16-1080>
48. Zhang, X., Jiang, Y., Wang, X., Hu, X., Sun, Y., Xie, P., Zhang, M.: Domain-specific NER via retrieving correlated samples. In: COLING (2022), <https://aclanthology.org/2022.coling-1.211>
49. Zhang, Y., Li, J., Song, Y., Zhang, C.: Encoding conversation context for neural keyphrase extraction from microblog posts. In: NAACL (Jun 2018). <https://doi.org/10.18653/v1/N18-1151>, <https://aclanthology.org/N18-1151>
50. Zheng, X., Mekala, D., Gupta, A., Shang, J.: News meets microblog: Hashtag annotation via retriever-generator. arXiv preprint (2021), <https://arxiv.org/abs/2104.08723>