

MULTPAX: Keyphrase Extraction using Language Models and Knowledge Graphs

Hamada M. Zahera^(✉) , Daniel Vollmers , Mohamed Ahmed Sherif , and
Axel-Cyrille Ngonga Ngomo 

DICE Group, Department of Computer Science,
Paderborn University, Germany
{hamada.zahera, daniel.vollmers, mohamed.sherif,
axel.ngonga}@uni-paderborn.de

Abstract. Keyphrase extraction aims to identify a small set of phrases that best describe the content of a text. The automatic generation of keyphrases has become essential for many natural language applications such as text categorisation, indexing, and summarization. In this paper, we propose MULTPAX, a multitask framework for extracting *present* and *absent* keyphrases using pretrained language models and knowledge graphs. In particular, our framework contains three components: first, MULTPAX identifies present keyphrases from an input document. Then, MULTPAX links with external knowledge graphs to get more relevant phrases. Finally, MULTPAX ranks the extracted phrases based on their semantic relatedness to the input document and return top- k phrases as a final output. We conducted several experiments on four benchmark datasets to evaluate the performance of MULTPAX against different state-of-the-art baselines. The evaluation results demonstrate that our approach significantly outperforms the state-of-the-art baselines, with a significance t-test $p < 0.041$. Our source code and datasets are public available at <https://github.com/dice-group/MultPAX>

Keywords: Present Keyphrase Extraction · Absent Keyphrase Generation · Knowledge Graph · Pretrained Language Models

1 Introduction

Keyphrase extraction is the process of extracting a small set of phrases that best describe a document. This process has been leveraged for several downstream applications, including text summarising, organising, and indexing [16]. In literature, the keyphrase extraction is divided into two sub-tasks: (i) detecting present keyphrases (PKE) that appear in a document, and (ii) generating absent keyphrases (AKG) that do not appear in the original document, but are essential for downstream applications (e.g., text summarization, indexing). Table 1 shows an example of extracting present and absent keyphrases from an input text.

Table 1: Example of present and absent keyphrase extraction from *Inspec* dataset. The predicted keyphrases are highlighted in blue, and the absent keyphrases are in red.

Input Text	<i>“This paper shows the importance that management plays in the protection of information and in the planning to handle a security breach when a theft of information happens. Recent thefts of information that have hit major companies have caused concern. These thefts were caused by companies’ inability to determine risks associated with the protection of their data and these companies lack of planning to properly manage a security breach when it occurs.”</i> quoted from [20]
Groundtruth Keyphrases	security breach, risk analysis, management issue, theft of information
Predicted Keyphrases	security breach, theft of information, security management, security risk, data management

Existing works mostly focus on extracting *present* keyphrases from input text, including supervised learning (e.g., sequence labelling [22]), and unsupervised learning (e.g., TextRank [17], YAKE [4]). By contrast, generating *absent* keyphrases (i.e., keyphrases that do not appear in a text) is a challenging task. According to the statistical study by [32], some benchmarking datasets (e.g., Inspec [9]) are missing up to 37.7% of *absent* keyphrases. To cope with this challenge, few studies have been proposed. For example, Meng et al. [15] employed a supervised sequence-to-sequence model with a *copy mechanism*, which allows copying important words directly from a source text, rather than decoding new words. However, this approach requires large-scale labelled data for training the model efficiently. In addition, the copy mechanism only generates one word at each time step and does not consider any dependencies between selected words [34]. Another line of work aims to utilise external sources of knowledge to generate absent keyphrases. For example, Shen et al. [25] constructs a phrase bank consisting of all keyphrases in a text corpus. The authors assumed that absent keyphrases in one document might be found in other relevant documents. However, this approach requires creating a domain-specific phrase bank to generate absent keyphrases.

In this paper, we aim not only to *extract* present keyphrases from an input document (in English), but also to *generate* absent keyphrases that are relevant and do not appear in the document. We reduce the effort required to develop a keyphrase model by employing precomputed resources. In particular, we use *pretrained language models* to extract present keyphrases and *knowledge graphs* (KGs) to generate absent keyphrases. For this purpose, we propose an unsupervised multitask framework (dubbed MULTPAX) with the following pipeline: i) We tokenize an input document into *n*-grams phrases and embed both (*document* and *n*-gram phrases) as low-dimensional vectors into one semantic space. Then,

we *extract* top- k phrases that are close to the document’s vector as candidates for present keyphrases. ii) We then *link* the extracted present keyphrases to find additional related terms (e.g., synonyms, hypernyms) from external KGs (e.g., DBPEDIA, BABELNET). For this purpose, we developed a new version of the MAG [18] framework, which is optimised for linking keywords and extracting related terms. iii) Finally, we *rank* all keyphrases (i.e., *present* and *absent*) based on their semantic similarity to the input document. The top- k phrases are returned as the final keyphrases output.

Additionally, we propose an improved metric for evaluating predicted keyphrases based on their *semantic-matching* with groundtruth keyphrases. Existing studies [13, 15, 33] consider *precision*, *recall*, and F_1 based on *exact-matching* between predicted and groundtruth keyphrases, which yields reasonable evaluation for present keyphrases that appear in text. However, in evaluating absent keyphrases, *exact-matching* demonstrated an inefficient assessment of words that are semantically similar but are literally different [21]. As an example, assume “*Cryptocurrency*” as a groundtruth keyphrase, and a keyphrase model was able to generate “*Bitcoin*” as a predicted keyphrase. In this case, *exact-matching* metric ignores the semantic relatedness between both words and considers them completely unrelated. By means of words embeddings, these words are similar and adjacent to each other in the embedding space. In this regard, we propose using an embedding-based F_1 score to evaluate keyphrases extraction in a more accurate semantic way.

To evaluate the performance of MULTPAX, we conducted several experiments on four benchmark datasets, where we study the performance of our system against different approaches. The evaluation results shows that our approach significantly outperforms the state-of-the-art baselines with a significance t-test $p < 0.041$ and F_1 -score up to 0.535. The main contributions in this paper can be summarised as follows:

- We propose an *unsupervised* multitask framework that not only extracts present keyphrases, but also generate absent ones.
- To the best of our knowledge, our approach is the first attempt that leverages existing *knowledge graphs* for keyphrases generation without the need to create keyphrases vocabularies or phrase banks.
- We introduce an *embedding-based* F_1 score that considers semantic similarity between generated and ground-truth keyphrases rather than the existing *exact-matching*.
- We carried out several experiments on four benchmark datasets. The evaluation results showed that our approach proved to be more accurate compared with state-of-the-art baselines.

2 Related Work

In this section, we give an overview of related approaches in unsupervised keyphrase extraction and absent keyphrase generation.

2.1 Unsupervised Keyphrase Extraction

Several approaches have recently been developed for extracting keyphrases in unsupervised setting without the need for annotated data. For example, statistical approaches such as TF-IDF and YAKE [4] compute statistical features (e.g., word frequencies and co-occurrences) to find important words as candidates for present keyphrases. Moreover, graph-based approaches like TextRank [17] construct a graph representation of text, where words are represented as nodes and their co-occurrences as edges. Thereafter, a node ranking algorithm (e.g., PageRank) is used to sort words, and return top- k words as candidate keyphrases. Bougouin et al. [3] proposed TopicRank, a graph-based approach similar to TextRank. In the first step, candidate phrases are clustered into topics and then ranked based on with their importance in the document.

Recent studies have demonstrated that embedding-based models can achieve significant results in extracting keyphrases. For example, EmbedRank [2] approach uses part-of-speech tags to extract potential keyphrases from an input document. Then, EmbedRank uses a pretrained embedding model to represent both phrases and an input document as low-dimensional vectors. Candidate keyphrases are then ranked based on their Cosine similarity scores to the document’s embedding vector. Although pretrained language models have shown promising performances for extracting present keyphrases, they have failed to generate absent keyphrases from their lexical corpus which are also related. Furthermore, Liang et al. [13] pointed out that embedding-based models ignore local information in a document. Accordingly, they developed a jointly-trained model to incorporate global and local context of a document. In the global view, their approach represented candidate keyphrases and a document as low-dimensional vectors into one semantic space. After that, the similarity between each candidate keyphrase and the document is computed. In terms of the local context, the authors built a graph structure based on the context of the document, where nodes represent phrases and edges represent similarities between them. Finally, the output keyphrases are ranked based on this global and local information.

2.2 Absent Keyphrase Extraction

Many previous approaches have relied on *sequence-to-sequence* models—with encoder-decoder architecture—to generate absent keyphrases [6]. By doing so, sequence-to-sequence models are able to decode not only keyphrases that appear in source text, but also those that may be absent, i.e., the ones that are not explicitly mentioned in the text. However, additional mechanisms need to be integrated to improve the generation of absent keyphrases. For example, Ye et al. [32] applied a *Graph Neural Network* (GNN) to capture knowledge from related references in scholarly publications. A *neural topic model* is employed in [29] to expand the context of the decoding component to generate more absent keyphrases.

It is noteworthy that Zhao et al. [33] achieved significant results in extracting keyphrases by dividing this task into two sub-tasks: present keyphrase extraction

and absent keyphrase generation. Furthermore, the authors proposed a multitask approach to *select*, *guide*, and *generate* keyphrases. In the *select* module, the authors used a BiLSTM to predict whether a sentence has a keyphrase or not. Then, a *guider* network is employed to utilize the attention information and memorize the predictions of the *selector*. Finally, this information is fed to a *generator* network to generate absent keyphrases by selecting words from both source text and a predefined vocabulary. In addition to these fully-supervised approaches, there are also some unsupervised methods that achieved promising results in generating keyphrases without the need for labelled data. Shen et al. [25] observed that many keyphrases absent from an input document appeared in other related documents. Therefore, they constructed a *phrase bank* of all keyphrases in a corpus. Then, they identified present keyphrases in relevant documents as candidates for absent keyphrases for the input document. In addition, they employed present keyphrases as *sliver labels* to train a sequence-to-sequence model. Finally, all keyphrases (both present and absent) were ranked based on their lexical and semantic similarity to an input document.

3 Our Approach

In this section, we present our approach for extracting *present* and *absent* keyphrases. Figure 1 depicts the architecture of our framework (MULTPAX), including three components: i) *present keyphrase extraction* (PKE), ii) *absent keyphrase generation* (AKG), and iii) *Keyphrases Semantic Matching*.

3.1 Problem Formulation

Let \mathcal{D} be an input document with $|S|$ sentences; each sentence $s \in S$ is a sequence of $|s|$ tokens $\mathcal{T} = \{t_1, t_2, \dots, t_{|s|}\}$. Our goal is to build a keyphrase model that not only extracts *present keyphrases* $\mathcal{Y}^p = \{y_1^p, y_2^p, \dots, y_{|\mathcal{Y}^p|}^p\}$ but also generates *absent keyphrases* $\mathcal{Y}^a = \{y_1^a, y_2^a, \dots, y_{|\mathcal{Y}^a|}^a\}$ that are relevant to \mathcal{D} by leveraging knowledge graphs such as DBPEDIA [1] and BABELNET [19].

Following previous works [8, 23], we divide the task of keyphrase extraction into two sub-tasks: *Present Keyphrase Extraction* (PKE) and *Absent Keyphrase Generation* (AKG). Furthermore, we define the computation of final keyphrases as a *Semantic Matching* task. First, we consider PKE as a *ranking* problem, where candidate phrases are extracted and then ranked based on their similarities to the input document (see Section 3.2). Second, we formulate AKE as a *linking* problem to infer relevant information from external knowledge graphs. For this task, we employ an unsupervised *entity-linker* [24] that maps a present keyphrase (\mathcal{Y}^p) to its corresponding entity in a knowledge graph (i.e., DBPEDIA, BABELNET) and then get relevant terms (e.g., from `dct:subject`, `gold:hypernym` properties) as candidates for absent keyphrases. Finally, all keyphrases $\mathcal{Y}^p \cup \mathcal{Y}^a$ are ranked based on their similarities to \mathcal{D} , top- k keyphrases are returned as the final output.

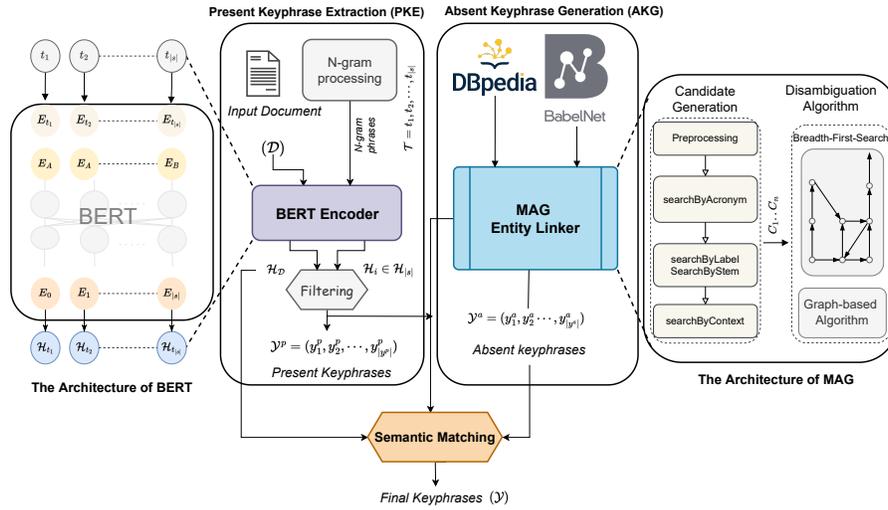


Fig. 1: The architecture of MULTPAX framework with three components: *present keyphrase extraction*, *absent keyphrase generation* and *semantic matching*.

3.2 Present Keyphrase Extraction (PKE)

We employ the BERT language model [7] to extract present keyphrases based on their semantic similarity to a document. The main steps are as follows: (1) We tokenize an input document \mathcal{D} into *n-gram phrases* and annotate each token with part-of-speech tags (e.g., ADJ: adjectives, NOUN: nouns, VERB: verbs). (2) Then, we remove stop words and keep noun phrases that consist of zero or more adjectives followed by one or multiple nouns [28]. (3) We employ the pretrained language model (BERT-Encoder) to encode candidate keyphrases as low-dimensional vectors together with the input document into one embedding space.

A special preprocessing is applied to the input text of the BERT-Encoder as follows: a [CLS] token is added at the beginning of each sentence, which is then used to obtain the contextualized embeddings vector of a sentence. An additional token [SEP] is inserted to mark the end of a sentence. Afterward, the input is tokenized by *WordPiece* tokenizer [26]; each token t_i is associated with three types of embeddings: *token embeddings* (E_{t_i}) which represents the vocabulary index of each token, *segmentation embeddings* that distinguishes between input sentences (E_A or E_B), and *position embeddings* (E_i) to indicate the position of each word. The output of the BERT-Encoder is the sentence’s representation matrix $\mathcal{H} = [h_0, h_1, \dots, h_{|s|}]$, where h_i denotes the embedding vector of token t_i . Formally, the embedding vector of a sentence s_j is

$$\mathcal{H}_j = \text{BERT-Encoder}(\{t_1, t_2, \dots, t_{|s|}\}) \quad (1)$$

Pooling is an essential operation for creating sentence and document embeddings [5]. It is commonly used to aggregate (e.g., mean, max) multiple

representations (e.g., sentences) into one embedding vector. To obtain the document’s embeddings $\mathcal{H}_{\mathcal{D}}$, we employ a `MaxPooling` layer on top of all sentences’ representations. Formally,

$$\mathcal{H}_{\mathcal{D}} = \text{MaxPooling}(\{\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_{|S|}\}). \quad (2)$$

Finally, we use Cosine distance to compute similarities between the embedding vectors of candidate keyphrases $\mathcal{H}_i \in H_{|S|}$ and the document embedding $\mathcal{H}_{\mathcal{D}}$. We select top- k keyphrases as candidates for present keyphrases.

3.3 Absent Keyphrase Generation (AKG)

To obtain absent keyphrases, we first link all present keyphrases \mathcal{Y}^p to a knowledge graph and get additional surface forms (i.e., strings that could be synonyms or alternative names). We consider DBPEDIA knowledge graph since it provides surface forms for a wide range of common entities. For *entity linking*, we follow a similar approach to the MAG framework [18].

MAG extracts entity links using two steps: *candidate generation* and *candidate disambiguation*. In the *candidate generation* step, MAG aims to find candidate links (C_1, \dots, C_n) for pre-marked entities in the search index [18]. To this end, MAG uses acronyms and labels in a knowledge graph to map premarked entity spans from the input text to candidate entities. Furthermore, MAG also relies on the Concise Bounded Description (CBD)¹ of the entities in a knowledge graph by comparing the context of the entity spans in the input document and the CBD of an entity in a knowledge graph [18]. We keep this candidate generation step from MAG and apply it to the extracted present keyphrases from the PKE component. In the *candidate disambiguation* step, MAG generates a local graph using a breadth-first-search method for all candidate entities on a knowledge graph. Then, MAG applies the HITS ranking algorithm [11] to jointly rank the candidate links for all entities in the local graph. HITS ranks the nodes in a directed graph based on incoming and outgoing edges. Authorities are seen as nodes, that carry important information, while hubs are nodes, that point to a large amount of authority nodes. So the authority score of a node n is calculated based on the hub score of the nodes, that have a directed edge to the node n , while the hub score of n is calculated based on the authority score of the nodes which are linked by n [11]. Formally, HITS calculates the authority score a_p for the node p as

$$a_p = \sum_{q:(q,p) \in G} h_q, \quad (3)$$

where h_q is the hub-score for the node q , given that a directed edge from node q to node p exists in the graph G . The hub-score h_p for a node p is calculated as

$$h_p = \sum_{q:(q,p) \in G} a_q, \quad (4)$$

¹ <https://www.w3.org/Submission/CBD/>

where a_q is the authority-score for a node q , which is linked by node p [11]. a_q and h_p are initialised randomly and updated iteratively until convergence.

In contrast to MAG, we not only link present keyphrases, but also extract related terms for each linked keyphrase from a knowledge graph. Furthermore, we extract top-ranked candidates for each entity and n nodes with the highest authority scores in the local graph, since their surface forms could be used as candidates for the absent keyphrases. In our approach, we use BABELNET to find *hypernyms* for the present keyphrases, in addition to the surface forms from DBPEDIA.

3.4 Keyphrases Semantic Matching

In the last component, we aim to identify top- k relevant keyphrases (*present* and *absent*), we set $k = \{5, 10, 20\}$ in our experiments. We regard this task as a semantic textual similarity [14]. To match similarities between document \mathcal{D} and candidate keyphrases, we embed them into one semantic space using a pretrained embedding model. Then we employ Cosine distance to find top- k nearest keyphrases (\mathcal{H}_i) to the document’s vector \mathcal{H}_D and return as final keyphrase predictions. Formally,

$$\text{Cos}(\mathcal{H}_i, \mathcal{H}_D) = \frac{\mathcal{H}_i \cdot \mathcal{H}_D}{\|\mathcal{H}_i\| \times \|\mathcal{H}_D\|} \quad (5)$$

where \mathcal{H}_i donates the embedding vector of candidate keyphrase (*present* y_i^p or *absent* y_i^a), and \mathcal{H}_D represents the embedding vector of the input document.

4 Experiments

We conducted our experiments to answer the following research questions:

- Q₁. *How efficient is our approach in extracting present keyphrases compared to the state-of-the-art approaches?*
- Q₂. *Are the existing exact-matching metrics (i.e., Precision, Recall and F₁-score) suitable for evaluating absent keyphrases?*
- Q₃. *To what extent does each component in our framework contribute to the overall performance?*

4.1 Experimental Setup

Datasets. In our experiments, we used four benchmark datasets of English documents, namely, *Inspec* [9], *SemEval2010* [10], *Krapivin* [12], and *NUS* [27]. Table 2 provides a statistical overview of each dataset, including the total number of documents (#Doc.), the number of documents in the evaluation (#Test), average keyphrases per document (Avg. KP) and the ratio of absent keyphrases in each dataset (Ratio%).

Table 2: Statistics about the datasets (#Doc: number of documents, #Test: size of test set, #Avg. KP: average keyphrase per document, #Ratio%: percentage of absent keyphrase per dataset).

Dataset	#Doc	#Test	Avg. KP	Ratio%
Inspec	2k	500	7.65	37.7%
Krapivin	2.3k	460	3.03	15.3%
SemEval2010	144	100	7.15	11.3%
NUS	211	211	2.71	17.8%

Baselines. We compared our approach with the following baselines for extracting keyphrases:

- **TextRank** [17] is an unsupervised approach that constructs a graph representation from a document, where nodes represent phrases and their edges are computed based on lexical similarities. Further, TextRank uses the *PageRank* algorithm to extract present keyphrases.
- **YAKE** [4] is a simple unsupervised method that extracts keywords automatically based on statistical features such as words co-occurrence and frequencies.
- **EmbedRank** [2] is an unsupervised method that employs words embeddings to identify relevant words to a document as candidate keyphrases. Furthermore, EmbedRank utilises the *Maximum Marginal Relevance* algorithm to increase the diversity of the extracted keyphrases.
- **Supervised-CopyRNN** [15] is a supervised baseline that trains a sequence-to-sequence model with a *copy mechanism* on KP20K dataset [15]. We used this approach as a baseline for present keyphrases extraction as well as absent keyphrase generation to compare the performance of copy mechanism.
- **AutoKeyGen** [25] is an unsupervised approach that constructs a *phrase bank* by combining keyphrases from all documents into a corpus. Then, AutoKeyGen considers lexical- and semantic-level similarities for selecting top candidate keyphrases (present and absent) for each input document.

Evaluation Metrics. We evaluated our approach using different metrics: *Precision*, *Recall*, and F_1 -score. *Precision* is computed as the number of correctly-matched keyphrases over all predicted keyphrases.

Given a list of predicted keyphrases $\mathcal{Y} = (y_1, \dots, y_{|\mathcal{Y}|})$, we select top- k ranked keyphrases $\mathcal{Y}_{:k} = (y_1, \dots, y_{\min(k, |\mathcal{Y}|)})$ and compare with top- k ranked keyphrases in the groundtruth set. We set $k = \{5, 10\}$ for present keyphrases and $k = \{10, 20\}$ for absent ones in our experiments. Following previous works [25, 31], we use the *Porter Stemmer* from the NLTK library² v3.7 to compute exact-matching between top- k predicted ($\mathcal{Y}_{:k}$) and groundtruth (\mathcal{Y}^{gold}) keyphrases. The precision of top- k predicted keyphrases is defined as

² <https://www.nltk.org/index.html>

$$P@k = \frac{|\mathcal{Y}_{:k} \cap \mathcal{Y}^{gold}|}{|\mathcal{Y}_{:k}|}. \quad (6)$$

Recall is calculated as how many correctly-matched keyphrases among all groundtruth keyphrases. Formally, the Recall is defined as

$$R@k = \frac{|\mathcal{Y}_{:k} \cap \mathcal{Y}^{gold}|}{|\mathcal{Y}_{:k}^{gold}|} \quad (7)$$

and $F_1@k$ score is defined as the harmonic mean of $P@k$ and $R@k$

$$F_1@k = 2 \times \frac{P@k \times R@k}{P@k + R@k} \quad (8)$$

Although *exact-matching* metric has been used widely in the literature [13], there is still a room for improvement regarding the absent keyphrases evaluation based on semantic similarity. Hence, we propose in Section 4.3 a semantic-based matching between keyphrase candidate and groundtruth to evaluate the performance of generated absent keyphrases.

Hyperparameters. We performed a grid search to optimise the hyperparameters of our approach. We found the following values yield the best F_1 scores. In the PKE component, we tokenized the input text into phrases of 2-4 grams. Further, we considered the top-10 ranked phrases as candidates for present keyphrases. The full setup of our experiments is available at the GitHub repository.³ For the baseline methods, the hyperparameters were set according to their original papers. In the MAG framework, we adapted the extraction of common entities to cover a larger set of entity types. In addition, we set the other hyperparameters values with the standard configuration⁴ of the MAG framework.

4.2 Present Keyphrase Evaluation (Q_1)

To answer Q_1 , we evaluated our approach (MULTPAX) vs. different baselines in extracting *present* keyphrases. As shown in Table 3, MULTPAX significantly outperforms all baselines by a large margin on most datasets with a significant *t-test* $p < 0.041$. This is due to, MULTPAX employs semantic similarity between candidate keyphrases and an input document using the state-of-the-art pretrained language model in semantic textual matching [30]. In contrast, CopyRNN [15] and AutoKeyGen [25] used sequence-to-sequence models to encode an input document as a low-dimensional vector and decode it back into a sequence of predicted keyphrases.

On the other hand, we find that YAKE does not perform well in detecting present keyphrases from short texts (e.g., papers’ abstracts). Since YAKE relies

³ <https://github.com/dice-group/MultPAX>

⁴ <https://github.com/dice-group/AGDISTIS/blob/master/src/main/resources/config/agdistis.properties>

Table 3: Evaluation results of *present* keyphrases prediction on *Inspec*, *SemEval2010*, *Krapivin*, and *NUS* datasets. $F_1@K$ scores are reported based on **exact-matching** between predicted and groundtruth keyphrases. Best results are reported in bold.

Model	Inspec		SemEval2010		Krapivin		NUS	
	$F_1@5$	$F_1@10$	$F_1@5$	$F_1@10$	$F_1@5$	$F_1@10$	$F_1@5$	$F_1@10$
TextRank	0.263	0.279	0.183	0.181	0.148	0.139	0.187	0.195
YAKE	0.027	0.038	0.050	0.242	0.013	0.020	0.013	0.020
EmbedRank	0.295	0.344	0.108	0.145	0.131	0.138	0.103	0.134
Supervised-CopyRNN	0.292	0.336	0.291	0.296	0.302	0.252	0.342	0.317
AutoKeyGen	0.303	0.345	0.187	0.240	0.171	0.155	0.218	0.233
MULTPAX	0.371	0.210	0.449	0.255	0.384	0.334	0.535	0.344

Table 4: *Absent* keyphrases evaluation (in terms of $R@10$, $R@20$). All results are reported based on **exact-matching** between predicted and groundtruth keyphrases, except the last row shows Recall results based on **semantic-matching**

Model	Inspec		SemEval2010		Krapivin		NUS	
	$R@10$	$R@20$	$R@10$	$R@20$	$R@10$	$R@20$	$R@10$	$R@20$
Supervised-CopyRNN	0.051	0.068	0.049	0.057	0.116	0.142	0.078	0.10
AutoKeyGen-Bank	0.015	0.017	0.007	0.009	0.031	0.041	0.021	0.026
AutoKeyGen-Full	0.017	0.021	0.010	0.011	0.033	0.054	0.024	0.032
MULTPAX _{exact-Matching}	0.079	0.080	–	–	–	–	0.017	0.017
MULTPAX _{semantic-Matching}	0.696	0.584	–	–	–	–	0.608	0.669

on statistical features such as words co-occurrence and frequencies, which are efficiently computed only in long texts (e.g., full papers or news). Remarkably, the embedding-based baseline (EmbedRank) achieves comparable results; however, it fails to generate absent keyphrases. In our approach, we extract present keyphrases from text using contextualised embeddings and semantic matching. These findings answer Q_1 ; by employing pretrained language models, we can not only efficiently identify present keyphrases from text without labelled data, but we also outperform the state-of-the-art approach (AutoKeyGen).

4.3 Absent Keyphrase Evaluation (Q_2)

We conduct further experiments to evaluate the performance of our approach against two baselines (namely, CopyRNN and AutoKeyGen) in generating absent keyphrases. Following previous work [25], we use Recall metric ($R@10$, $R@20$) based on *exact-matching* for the performance evaluation as shown in Table 4. Since we used the same experimental setup of CopyRNN and AutoKeyGen approaches, we obtained the evaluation results from their papers [15, 25].

Regarding Q_2 , we can clearly see that all approaches achieve poor performances when considering exact-matching between predicted and groundtruth keyphrases. For example, if two keyphrases are semantically similar, e.g., “*disaster relief organization*” and “*crisis responses institute*”, these keyphrases will not be considered as a *match* using the existing metrics. Hence, we found that such metrics are unsuitable for evaluating absent keyphrases. We propose an improved evaluation metric based on *semantic-matching*. Formally, let \mathcal{Y}^a be predicted keyphrases; \mathcal{Y}^{gold} is groundtruth keyphrases. We first embed each keyphrase in \mathcal{Y}^a and \mathcal{Y}^{gold} . Then, we use Cosine distance to compute similarities between the embedding of each keyphrase in \mathcal{Y} and \mathcal{Y}^{gold} . We set a threshold (> 0.5) for similarities scores to consider semantic-matching between \mathcal{Y} and \mathcal{Y}^{gold} . The two last rows in Table 4 show the evaluation results of R@10 and R@20 based on semantic-matching compared to exact-matching in absent keyphrase extraction.

The AutoKeyGen baseline demonstrates competitive performance in generating absent keyphrases on the NUS dataset. However, the generated keyphrases by AutoKeyGen are limited to the ones in the phrase bank of each dataset. In contrast, our approach leverages public knowledge graphs (such as DBPEDIA and BABELNET) to obtain relevant phrases as candidates for absent keyphrases.

Limitation of our work. In our experiment, we used the MAG framework to connect present keyphrases to DBPEDIA knowledge graph (see Section 3.3). In the SemEval2010 and Krapivin datasets, we were unable to link present keyphrases, due to the lack of coverage for these keyphrases in the DBPEDIA knowledge graph. That is the reason for the missing values shown in the last two rows of Table 4 for these datasets. In our future work, we plan to integrate other knowledge graphs (e.g., YAGO and WIKIDATA) to extend the coverage of entity linking in the MAG framework.

4.4 Ablation Study (Q_3)

To answer Q_3 , we analysed the impact of each component of our framework on the overall performance. For this purpose, we set up four variants of our framework. The first variant MULTPAX-PKE was dedicated for only extracting present keyphrases, i.e., no absent keyphrase generation and thus no linking with knowledge graphs. We also created two variants of MULTPAX with the purpose of testing the generating of absent keyphrases, namely MULTPAX-AKE_{DBpedia} and MULTPAX-AKE_{BabelNet}. Furthermore, we configured the MAG framework to link present keyphrases only with DBPEDIA in case of MULTPAX-AKE_{DBpedia}, and only with BABELNET for MULTPAX-AKE_{BabelNet}. Finally, we benchmarked the entire framework MULTPAX_{Full} as our fourth variant.

Table 5 reports the evaluation results of each component in terms of *semantic-matching* $F_1@5$, and $F_1@10$ on the *Inspec* dataset, since it contains the highest ratio of absent keyphrases among the benchmark datasets. We can see that the performance of MULTPAX-PKE is improved when linking with knowledge graphs, e.g., MULTPAX-AKE_{DBpedia} outperforms MULTPAX-PKE by +0.41 in $F_1@10$. In addition, we noticed that our approach could retrieve more terms from

Table 5: Ablation Study of MULTPAX framework on *Inspec* dataset. $F_1@K$ scores are reported based on **semantic-matching** between predicted and groundtruth keyphrases.

MULTPAX-variant	$F_1@5$	$F_1@10$
MULTPAX-PKE	0.892	0.686
MULTPAX-AKE _{BabelNet}	0.907	0.701
MULTPAX-AKE _{DBpedia}	0.911	0.727
MULTPAX_{Full}	0.911	0.763

DBPEDIA than BABELNET, since DBPEDIA contains more semantic ontologies (approximately 3.5 millions instances) extracted from Wikipedia information boxes. Finally, our MULTPAX_{Full} showed an improved performance with F_1 scores (0.911 in $F_1@5$, 0.763 in $F_1@10$) when incorporating both knowledge graphs (i.e., DBPEDIA and BABELNET) compared with individual variants. These findings conclude that each component of MULTPAX contributes to the overall performance of our framework and answers our last research question Q_3 .

5 Conclusion

This paper presents MULTPAX framework, a multitask approach for extracting present and absent keyphrases, including three components: i) Present Keyphrase Extraction, ii) Absent Keyphrases Generation, and iii) Keyphrases Semantic Matching. In our approach, we employ a pretrained language model (BERT) and knowledge graphs (DBPEDIA and BABELNET) in keyphrase extraction. Our experiments showed that pretrained language models are capable of efficiently extracting present keyphrases. Furthermore, knowledge graphs proved to be valuable resources for generating keyphrases that are absent, especially in short text. In our future work, we plan to apply a bootstrapped approach for keyphrase extraction from DBPEDIA abstracts to find more relevant terms. In particular, we intend to apply MULTPAX recursively on the abstracts of DBPEDIA entities. In addition, we will experiment with other knowledge graphs (e.g., YAGO and WIKIDATA) to extend the coverage of entity link in the MAG framework.

Supplemental Material Statement. We implemented our framework in Python 3.7, the source code and how-to-run instructions can be found at the GitHub repository.⁵ Furthermore, we used the benchmarking datasets available in the Dropbox drive.⁶ For the baseline models, we used OpenNMT library⁷, which enables us to benchmark different state-of-the-art baselines in our experiments. In addition, we used the pretrained embedding of BERT model,

⁵ <https://github.com/dice-group/MultPAX>

⁶ <https://www.dropbox.com/s/aluvkblymjs7i3r/MULTPAX-Datasets.zip?dl=0>

⁷ <https://github.com/memray/OpenNMT-kpg-release>

namely all-MiniLM-L6-v2 with 384 embedding dimension from the `huggingface library`⁸ v4.16. For the hardware requirements, we used a computing server with 256 GB memory and Xeon(R) CPU E5-2630 v4 with 2.20 GHz to run our experiments.

Acknowledgments This work has been supported by the German Federal Ministry for Economic Affairs and Climate Action (BMWK) within the projects RAKI (grant no 01MD19012B) and SPEAKER (grant no 01MK20011U) as well as by the German Federal Ministry of Education and Research (BMBF) within the projects COLIDE (grant no 01I521005D) and EML4U (grant no 01IS19080B). We are also grateful to Diego Moussallem for the valuable discussion on earlier drafts and Pamela Heidi Douglas for editing the manuscript.

References

1. Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z.: Dbpedia: A nucleus for a web of open data. In: *The semantic web*, pp. 722–735, Springer (2007)
2. Bennani-Smires, K., Musat, C., Hossmann, A., Baeriswyl, M., Jaggi, M.: Simple unsupervised keyphrase extraction using sentence embeddings. In: *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pp. 221–229 (2018)
3. Bougouin, A., Boudin, F., Daille, B.: Topicrank: Graph-based topic ranking for keyphrase extraction. In: *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, pp. 543–551 (2013)
4. Campos, R., Mangaravite, V., Pasquali, A., Jorge, A., Nunes, C., Jatowt, A.: Yake! keyword extraction from single documents using multiple local features. *Information Sciences* **509**, 257–289 (2020)
5. Chen, Q., Ling, Z.H., Zhu, X.: Enhancing sentence embedding with generalized pooling. In: *Proceedings of the 27th International Conference on Computational Linguistics*, pp. 1815–1826 (2018)
6. Chen, W., Gao, Y., Zhang, J., King, I., Lyu, M.R.: Title-guided encoding for keyphrase generation. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 6268–6275 (2019)
7. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018)
8. Gollapalli, S.D., Li, X.L., Yang, P.: Incorporating expert knowledge into keyphrase extraction. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31 (2017)
9. Hulth, A.: Improved automatic keyword extraction given more linguistic knowledge. In: *Proceedings of the 2003 conference on Empirical methods in natural language processing*, pp. 216–223 (2003)
10. Kim, S.N., Medelyan, O., Kan, M.Y., Baldwin, T.: Semeval-2010 task 5: Automatic keyphrase extraction from scientific articles. In: *Proceedings of the 5th International Workshop on Semantic Evaluation*, pp. 21–26 (2010)

⁸ <https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

11. Kleinberg, J.M.: Authoritative sources in a hyperlinked environment. *Journal of the ACM (JACM)* **46**(5), 604–632 (1999)
12. Krapivin, M., Autaeu, A., Marchese, M.: Large dataset for keyphrases extraction (2009)
13. Liang, X., Wu, S., Li, M., Li, Z.: Unsupervised keyphrase extraction by jointly modeling local and global context. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 155–164 (2021)
14. Majumder, G., Pakray, P., Gelbukh, A., Pinto, D.: Semantic textual similarity methods, tools, and applications: A survey. *Computación y Sistemas* **20**(4), 647–665 (2016)
15. Meng, R., Zhao, S., Han, S., He, D., Brusilovsky, P., Chi, Y.: Deep keyphrase generation. In: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 582–592 (2017)
16. Merrouni, Z.A., Frikh, B., Ouhbi, B.: Automatic keyphrase extraction: a survey and trends. *Journal of Intelligent Information Systems* **54**(2), 391–424 (2020)
17. Mihalcea, R., Tarau, P.: Textrank: Bringing order into text. In: *Proceedings of the 2004 conference on empirical methods in natural language processing*, pp. 404–411 (2004)
18. Moussallem, D., Usbeck, R., Röder, M., Ngonga Ngomo, A.C.: Mag: A multilingual, knowledge-base agnostic and deterministic entity linking approach. In: *K-CAP 2017: Knowledge Capture Conference*, p. 8, ACM (2017)
19. Navigli, R., Ponzetto, S.P.: Babelnet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artificial intelligence* **193**, 217–250 (2012)
20. Polstra III, R.M.: A case study on how to manage the theft of information. In: *Proceedings of the 2nd annual conference on Information security curriculum development*, pp. 135–138 (2005)
21. Ray Chowdhury, J., Caragea, C., Caragea, D.: Keyphrase extraction from disaster-related tweets. In: *The world wide web conference*, pp. 1555–1566 (2019)
22. Sahrawat, D., Mahata, D., Zhang, H., Kulkarni, M., Sharma, A., Gosangi, R., Stent, A., Kumar, Y., Shah, R.R., Zimmermann, R.: Keyphrase extraction as sequence labeling using contextualized embeddings. In: *European Conference on Information Retrieval*, pp. 328–335, Springer (2020)
23. Sahrawat, D., Mahata, D., Zhang, H., Kulkarni, M., Sharma, A., Gosangi, R., Stent, A., Kumar, Y., Shah, R.R., Zimmermann, R.: Keyphrase extraction as sequence labeling using contextualized embeddings. In: *European Conference on Information Retrieval*, pp. 328–335, Springer (2020)
24. Shen, W., Wang, J., Han, J.: Entity linking with a knowledge base: Issues, techniques, and solutions. *IEEE Transactions on Knowledge and Data Engineering* **27**(2), 443–460 (2014)
25. Shen, X., Wang, Y., Meng, R., Shang, J.: Unsupervised deep keyphrase generation. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, pp. 11303–11311 (2022)
26. Song, X., Salcianu, A., Song, Y., Dopson, D., Zhou, D.: Fast wordpiece tokenization. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 2089–2103 (2021)
27. Vijayakumar, A.K., Cogswell, M., Selvaraju, R.R., Sun, Q., Lee, S., Crandall, D., Batra, D.: Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424* (2016)
28. Wan, X., Xiao, J.: Single document keyphrase extraction using neighborhood knowledge. In: *AAAI*, vol. 8, pp. 855–860 (2008)

29. Wang, Y., Li, J., Chan, H.P., King, I., Lyu, M.R., Shi, S.: Topic-aware neural keyphrase generation for social media language. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 2516–2526 (2019)
30. Xia, T., Wang, Y., Tian, Y., Chang, Y.: Using prior knowledge to guide bert’s attention in semantic textual matching tasks. In: Proceedings of the Web Conference 2021, pp. 2466–2475 (2021)
31. Ye, H., Wang, L.: Semi-supervised learning for neural keyphrase generation. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pp. 4142–4153 (2018)
32. Ye, J., Cai, R., Gui, T., Zhang, Q.: Heterogeneous graph neural networks for keyphrase generation. arXiv preprint arXiv:2109.04703 (2021)
33. Zhao, J., Bao, J., Wang, Y., Wu, Y., He, X., Zhou, B.: Sgg: Learning to select, guide, and generate for keyphrase generation. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 5717–5726 (2021)
34. Zhao, Y., Song, J., Feng, H., Zhuang, F., Li, Q., Wang, X., Liu, J.: Deep keyphrase completion. arXiv preprint arXiv:2111.01910 (2021)