

KPDROP: Improving Absent Keyphrase Generation

Jishnu Ray Chowdhury[♣] Seoyeon Park^{♣*} Tuhin Kundu^{♣*†} Cornelia Caragea[♣]

[♣] Computer Science, University of Illinois at Chicago [♣] Amazon

{jraych2, spark313, cornelia}@uic.edu

tuhinkundu@outlook.com

Abstract

Keyphrase generation is the task of generating phrases (keyphrases) that summarize the main topics of a given document. Keyphrases can be either present or absent from the given document. While the extraction of present keyphrases has received much attention in the past, only recently a stronger focus has been placed on the generation of absent keyphrases. However, generating absent keyphrases is challenging; even the best methods show only a modest degree of success. In this paper, we propose a model-agnostic approach called keyphrase dropout (or KPDROP) to improve *absent keyphrase generation*. In this approach, we randomly drop present keyphrases from the document and turn them into artificial absent keyphrases during training. We test our approach extensively and show that it consistently improves the absent performance of strong baselines in both supervised and resource-constrained semi-supervised settings¹.

1 Introduction

Keyphrase generation (KG) is the task of producing a set of phrases that best summarize a document. It can be leveraged for various applications such as text summarization (Zhang et al., 2017), recommendation (Bai et al., 2018), and opinion mining (Meng et al., 2012). Accurate identification of keyphrases especially on scientific papers can improve efficiency in paper indexing and paper retrieval (Chen et al., 2019a; Boudin et al., 2020).

Keyphrases can be divided into two types: (1) *present* keyphrases and (2) *absent* keyphrases. Present keyphrases appear verbatim in the document, whereas absent keyphrases are topically-relevant but missing from the document. Many prior works (Hasan and Ng, 2014; Ünlü and Çetin,

2019) focused on present keyphrase *extraction* exclusively. As such, they are not suitable for generating any absent keyphrases. To overcome this limitation, recent approaches (Meng et al., 2017; Yuan et al., 2020; Chen et al., 2020; Ye et al., 2021b) use sequence-to-sequence (seq2seq) models to generate both present and absent keyphrases.

As shown by Boudin and Gallina (2021), absent keyphrases that are substantially different from the present ones can significantly improve the effectiveness of document-retrieval. Boudin and Gallina (2021) suggested that absent keyphrases can improve document-retrieval by expanding the query terms to alleviate the vocabulary mismatch problem between the query terms and relevant documents (Furnas et al., 1987). Thus, there is a strong motivation for improving absent keyphrase performance. However, generating absent keyphrases can be very challenging. Even the best keyphrase generation models (Chen et al., 2020; Ye et al., 2021a,b) still only achieve a modest performance in absent keyphrases.

In this work, we propose *keyphrase dropout* or KPDROP as a simple and effective technique to improve the performance of *absent* keyphrases. Unlike the traditional dropout method that focuses on dropping neurons (Srivastava et al., 2014), we propose a novel dropout method where, instead of neurons, we randomly drop entire phrases (specifically, present gold keyphrases) from a given document during training. Thus, the dropped present keyphrases are turned into (artificial) absent keyphrases. As a result, KPDROP has the following two effects:

1. The model is forced to deeply utilize the context information to infer dropped keyphrases that could have been otherwise simply extracted from the text. Thereby, the capability to infer missing keyphrases in general (including naturally missing ones) is increased.

* Equal contribution

† Work done at the University of Illinois at Chicago before Amazon

¹Code: <https://github.com/JRC1995/KPDrop>

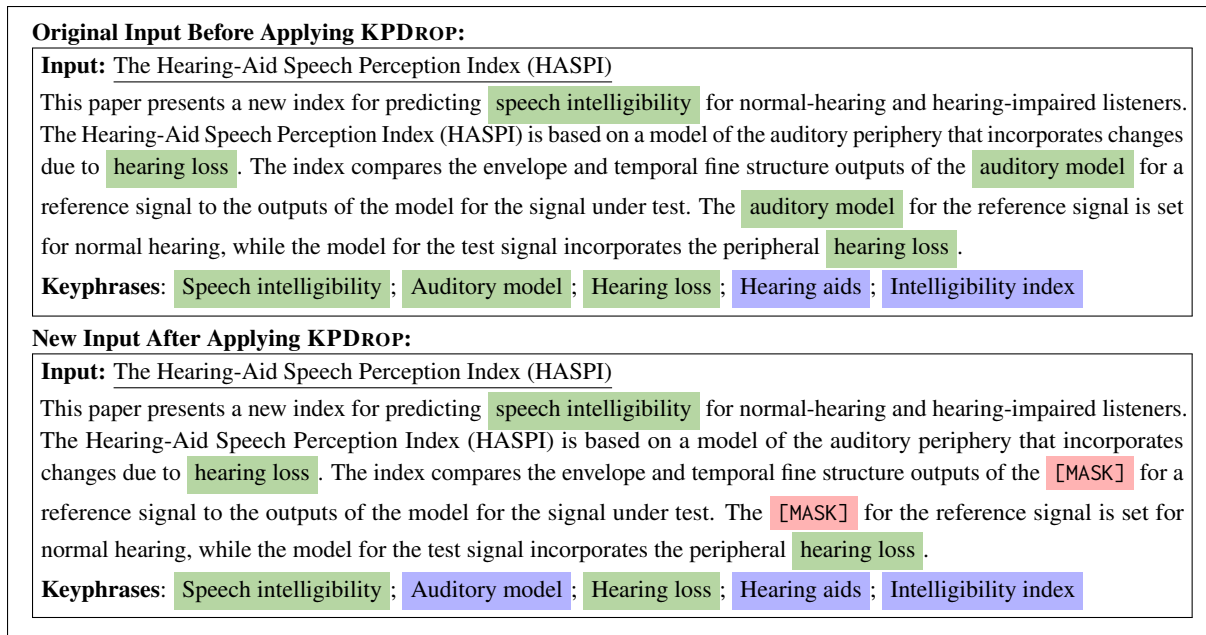


Figure 1: An example of how the input document (Kates and Arehart, 2014), and the keyphrases change after applying KPDR0P. Here, KPDR0P drops the present keyphrase “auditory model”. Green highlighting indicates present keyphrases, blue highlighting indicates absent keyphrases, and red highlighting indicates mask tokens.

2. The method stands as a dynamic data augmentation strategy. KPDR0P can be used to modify a given set of training documents by randomly dropping some present keyphrases from the documents turning them into artificial absent keyphrases. As such, new data can be created from the originals just like a standard data augmentation technique. Moreover, the process can be dynamically done in real-time during training. Thus, in different epochs different present keyphrases may be dropped for the same document yielding higher diversity in the data for model training.

We apply KPDR0P into three distinct neural models representing three distinct paradigms: one2many (Yuan et al., 2020), one2one (Meng et al., 2017; Huang et al., 2021), and one2set (Ye et al., 2021b). We observe consistent and substantial improvement in absent keyphrase performance in supervised settings on five standard datasets used for KG evaluation, with little to no drops in present keyphrase performance (and in fact often yielding improved present performance). Additionally, KPDR0P can be used to create synthetic absent keyphrases for unlabelled data to be used for self-supervised pre-training. We demonstrate that such pre-training augmented with KPDR0P brings substantial improvement when fine-tuned in low-resource labelled data.

2 Methodology of KPDR0P

In this section, we formally describe our approach of Keyphrase Dropout (KPDR0P). In keyphrase generation (KG), we have a document X as an input sequence of tokens, and a set of n keyphrases, $Y = \{y_1, y_2, \dots, y_n\}$, as the target output. Within Y , there is a subset of s ($0 \leq s \leq n$) present keyphrases as $P = \{p_1, p_2, p_3, \dots, p_s\}$ ($P \subseteq Y$) and a subset of t ($0 \leq t \leq n$) absent keyphrases as $A = \{a_1, a_2, a_3, \dots, a_t\}$ ($A \subseteq Y$). It also holds true that $A \cap P = \emptyset$ and $s + t = n$. Similar to Meng et al. (2017), we separate absent and present keyphrases by checking if the stemmed version of a keyphrase appears in the stemmed version of the input document (in which case it is present) or does not appear (in which case it is absent).

When applying KPDR0P, any present keyphrase $p_k \in P$ can be randomly dropped with a probability of r (sampled from binomial distribution), where r is the dropout rate (or KPDR0P rate). Let O ($O \subseteq P$) be the subset of present keyphrases that are randomly dropped during the application of KPDR0P at some training epoch. For every keyphrase $p_k \in O$, we remove *any and all* substrings from X that, when stemmed, match the stemmed version of the phrase p_k and replace each removed substring with a special mask token [MASK]. We call the modified version of X as X^{new} . The sets of present and

absent keyphrases are also modified. The new set of present keyphrases becomes $P := P \setminus O$ and the new set of absent keyphrases becomes $A := A \cup O$. Thus, the keyphrases in O , which were originally present, become artificially absent after applying KPDROP. The set Y becomes Y^{new} with the new versions of P and A . An example of applying KPDROP is shown in Figure 1.

Note that although as a set Y^{new} is the same as Y , we use Y^{new} to convey that the sets of present and absent keyphrases have changed.

2.1 Two Strategies of Applying KPDROP

KPDROP can be applied to a given mini-batch of examples in at least two distinct ways which act as data augmentation (noising) strategies.

KPDROP-R: For the first strategy that we refer to as KPDROP-REPLACE or KPDROP-R, we can think of KPDROP as a dropout technique and like other applications of dropout we can choose to *replace* the original examples with their corresponding keyphrase-dropped versions. In other words, each original sample (X, Y) in a mini-batch B is replaced with (X^{new}, Y^{new}) obtained after applying KPDROP. More formally, if initially we had a mini-batch set B , after applying KPDROP-R, we get a new batch B_{KPDROP} as $B_{KPDROP} = \{(X^{new}, Y^{new}) | (X, Y) \in B\}$. Note that at any given training iteration, we create a *single* KPDropped counterpart (X^{new}, Y^{new}) for each sample (X, Y) in the mini-batch.

KPDROP-A: For the second strategy that we refer to as KPDROP-APPEND or KPDROP-A, we can think of KPDROP as closer to a data augmentation technique where instead of replacing the original examples we *augment* the original batch with the new keyphrase dropped versions of the originals. In other words, for each original sample (X, Y) in a mini-batch, we add (X^{new}, Y^{new}) obtained after applying KPDROP. More formally, starting with the same mini-batch set B as before, after applying KPDROP-A, we get a new batch $B_{KPDROP-A}$ as $B_{KPDROP-A} = B \cup B_{KPDROP}$.

Given that KPDROP-R can increase the number of absent keyphrases per sample during training, we hypothesize that the model can learn to be more biased towards generating absent keyphrases. This can help to increase absent keyphrase performance at the cost of present keyphrase performance because the technique will be dropping present

keyphrases. On the other hand, KPDROP-A, instead of replacing the original data, it can offer extra samples per batch that have additional artificially absent keyphrases. Thus, KPDROP-A should be still able to improve absent keyphrase performance while also maintaining the original data with its original present keyphrases. Thus, KPDROP-A will offer the underlying model with the same opportunity to learn present keyphrases as would the model without KPDROP-A. Intuitively, this can help improve absent performance without a substantial cost to present performance.

2.2 KPDROP Features and Connections to Works from Other Areas

In this section, we highlight some notable features of KPDROP and draw connections to some relevant works outside the area of keyphrase generation.

KPDROP and Masked Language Modeling: Superficially, KPDROP is similar to a standard masked language modeling (MLM) (Devlin et al., 2019; Raffel et al., 2020) strategy insofar that both involve the reconstruction of some masked out text-spans. However, there are several crucial differences between KPDROP and vanilla MLM. First, KPDROP masks only *present keyphrases*, not random subspans or phrases. Second, KPDROP masks *all instances* of the selected present keyphrases in the document to make them truly absent (see Figure 1). In contrast, MLM does not mask *all instances* of the masked token or phrase. Third, masking (replacing a subspan with a mask token) is only an engineering choice for KPDROP; not an essential ingredient. We can simply drop the selected present keyphrases without replacing them with a mask token. In KPDROP, predictions of artificial absent keyphrases (masked present keyphrases) are not *explicitly* associated with any mask position.

KPDROP as Structured Dropout: KPDROP is a unique kind of *structured dropout* where all instances of some randomly chosen present keyphrases are dropped. There are other examples of structured dropout in prior works. For example, Iyyer et al. (2015) proposed word dropout (for general NLP tasks) where whole word embeddings are dropped instead of just random neurons. Huang et al. (2016) and Fan et al. (2020) used another form of structured dropout that stochastically drops whole layers for general computer vision tasks and for general NLP tasks respectively.

Model-Agnosticism of KPDRÖP: One important feature of KPDRÖP is that it is model-agnostic. KPDRÖP only requires changing the input and output without making internal architectural changes. As such, it is compatible with any model that is suitable for KG. KPDRÖP can also be easily stacked with other techniques that help absent keyphrase generation. Later (in §5.1), we demonstrate that KPDRÖP works just as well for very different architectures: (1) an RNN-based model trained in one2many settings (Yuan et al., 2020) (2) an RNN-based model trained in one2one settings (Meng et al., 2017; Huang et al., 2021), and (3) a Transformer-based model trained in one2set settings (Ye et al., 2021b).

3 Evaluations

We use the following evaluation metrics:

F₁@M: F₁@M is a standard metric (Yuan et al., 2020; Chen et al., 2019b; Chan et al., 2019; Chen et al., 2020; Ye et al., 2021b) used to evaluate the performance of keyphrase generation. This is an F₁ based metric where *all* the predictions by a given model are considered for evaluation.

F₁@5: F₁@5 is similar to F₁@M but only the top 5 predicted keyphrases are used for evaluation (if the total number of predictions are less than 5, all the available predictions are used). This metric is used in several prior works (Meng et al., 2017; Yuan et al., 2020).² This metric is useful in settings where the model is used to overgenerate keyphrases (for example, by using beam search). Overgeneration can lead to lower precision when all the predictions are kept (thus, low F₁@M). So, often in such settings, it is more useful to check the performance of the model when some simple truncation policy is used, for example, selecting some top *k* (e.g., top 5 in F₁@5) predictions.

R@10 and R@50: R@10 and R@50 represents macro recall of the top 10 predictions and the top 50 predictions, respectively. In some applications (such as retrieval) high recall can be sometimes more useful than precision. R@10 and R@50 are used in prior works (Meng et al., 2017; Chen et al., 2018; Liu et al., 2020) to measure absent keyphrase

²Note that the F₁@5 metric as used by us and as introduced by Meng et al. (2017) for keyphrase generation is different from the F₁@5 metric as introduced by Chan et al. (2019) and used in some other works (Chen et al., 2020; Ye et al., 2021b). We report results with Chan et al. (2019)’s formulation of F₁@5 in Appendix A.2

performance after overgenerating them using beam search with high beam size.

Like prior works, we use macro-F₁ and macro-recall for all the above metrics.

4 Baselines

We use the following baselines with our KPDRÖP-R and KPDRÖP-A approaches:

GRU One2Many: GRU One2Many (also known as CatSeq) represents a simple seq2seq baseline based on GRUs. It takes a document input and generates a concatenated sequence of keyphrases similar to Yuan et al. (2020). For this baseline, we concatenate the ground truth keyphrases based on the best performing ordering procedure (Meng et al., 2021): we first concatenate the present keyphrases according to the order of their first appearance in text and then we append the absent keyphrases in their natural order. However, when using KPDRÖP, we start with the ordering as mentioned but then shift the dropped present keyphrases (artificial absent keyphrases) after the fully present keyphrases but keep them before the naturally absent keyphrases. We maintain the internal order of the dropped present (artificially absent) keyphrases. We use “;” as the delimiter for separating keyphrases.

GRU One2One: GRU One2One (also known as CopyRNN) is another seq2seq model based on GRUs. However, unlike the One2Many model, it can predict only one keyphrase per generated sequence. It can be still used to generate multiple keyphrases by using beam search and preserving multiple beams of sequences (each representing a keyphrase). The One2One approach was first introduced by Meng et al. (2017) where the original training data was divided such that each input was associated with only one ground truth keyphrase. However, instead of that, we use the more efficient training approach of using reset states as in Huang et al. (2021). That is, we simply train the one2one model similar to one2many models using teacher forcing but we “reset” the hidden state when generating the first word of any ground truth keyphrase. Resetting (Huang et al., 2021) corresponds to using the initial RNN hidden state and the first special input token indicating start of sequence and thereby removing dependencies from previous generations.

Transformer One2Set: Transformer One2Set (also known as SetTrans) is a Transformer-based

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5
GRU One2Many										
Greedy	1.4 ₂	1.4 ₂	2.9 ₆	2.9 ₆	3.4 ₄	3.4 ₄	2.4 ₂	2.4 ₂	3.3 ₃	3.3 ₃
Greedy+KPD-R	1.2 ₁	1.2 ₁	3.2 ₇	3.2 ₇	5.1₂	5.1₂	2.7₂	2.7₂	4.0₁	4.0₁
Greedy+KPD-A	1.5₁	1.5₁	3.5₆	3.5₆	4.6 ₃	4.6 ₃	2.6 ₄	2.6 ₄	3.9 ₃	3.9 ₃
Beam	2.8 ₁	2.9₁	5.5 ₃	5.5 ₂	7.2 ₃	7.3 ₄	3.8 ₂	3.8 ₂	5.8 ₀	5.8 ₁
Beam+KPD-R	2.8 ₃	2.9₃	6.5 ₂	6.5₂	7.3₃	7.4₃	5.1₂	4.8 ₂	6.1 ₀	6.3 ₀
Beam+KPD-A	3.0₁	2.9₁	6.6₁₀	6.5₉	7.3₃	7.3 ₄	5.0 ₃	5.0₄	6.4₀	6.5₀
GRU One2One										
Beam	0.6 ₀	2.8 ₁	1.5 ₁	5.7 ₇	1.2 ₀	5.9 ₃	1.3 ₀	3.8 ₁	1.0 ₀	6.2 ₀
Beam+KPD-R	0.7₀	2.9₃	1.5 ₀	7.5₃	1.3₀	7.8₂	1.4 ₀	4.9₅	1.0 ₀	6.5 ₁
Beam+KPD-A	0.7₀	2.7 ₁	1.6₀	6.7 ₂	1.3₀	6.5 ₃	1.5₁	4.1 ₃	1.1₀	6.8₀
Transformer One2Set										
Greedy	2.8 ₄	2.8 ₄	6.4 ₈	6.4 ₈	6.8 ₂	6.8 ₂	3.5 ₁	3.5 ₁	5.6 ₀	5.5 ₀
Greedy+KPD-R	2.9 ₁	2.8 ₁	6.9 ₈	6.9 ₇	8.4₈	8.4₈	4.6 ₃	4.6₄	6.4 ₁	6.3 ₁
Greedy+KPD-A	3.2₂	3.2₂	7.4₉	7.4₉	7.2 ₇	7.2 ₇	4.7₁	4.6₁	6.6₁	6.5₁
Beam	0.4 ₀	3.3 ₂	0.9 ₀	7.0 ₅	0.8 ₀	6.7 ₅	0.8 ₀	4.7 ₄	0.6₀	5.8 ₀
Beam+KPD-R	0.4 ₀	2.4 ₁	1.0₀	7.2 ₅	0.9₀	7.3 ₂	0.9 ₀	5.2 ₄	0.6₀	6.1 ₀
Beam+KPD-A	0.5₀	3.6₂	1.0₀	7.9₄	0.9₀	7.8₂	1.0₀	5.3₄	0.6₀	6.7₀

Table 1: Absent keyphrase performance (F1) for different models. KPD represents KPDrop. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1). We bold the best scores per block.

model trained in a new One2Set paradigm as introduced by Ye et al. (2021b). In this method, a fixed number of preset control codes are used to generate all present and absent keyphrases in parallel independent of each other. Moreover, during training target keyphrases are matched with the predicted ones using Hungarian algorithm (Kuhn, 1955) so that the training is not sensitive to the order of the target keyphrases. Note that Transformer One2Set uses specialized control codes for present and absent keyphrase generation separately. Thus, when using KPDROP, we make sure to use the control codes for absent keyphrases to generate artificial absent keyphrases (dropped present keyphrases) as well.

5 Supervised Experiments

In the supervised setting, we explore the effects of applying both KPDROP-R and KPDROP-A on the baselines that we discussed in §4. Note that One2Many models (Yuan et al., 2020; Meng et al., 2021) had been explored in both greedy search based generations where only a single concatenated sequence of keyphrases is generated and also in beam search based generations where multiple beams of concatenated sequence of keyphrases are generated. We too explore both greedy and beam search for the One2Many models. In One2One models, greedy search is ineffective because it

only generates a single keyphrase. Thus, following Meng et al. (2017), we only use beam search to overgenerate multiple beams of keyphrases for One2One models. For One2Set models, only greedy search was explored (Ye et al., 2021b). Here, in addition to greedy search, we also explore using beam search for each control code in One2Set models. For all models, we investigate the effect of applying KPDROP in all of their applicable decoding settings (be it beam search or greedy search). Following prior works (Meng et al., 2017; Chen et al., 2018), we mainly use beam search to demonstrate recall performance (Table 2) for absent keyphrase performance. The recall performance stands out best when using beam search. Greedy search can be more conservative with generation; thus, has limited recall. We evaluate our three baselines on five scientific datasets: KP20k (Meng et al., 2017), Inspec (Hulth, 2003), Krapivin (Krapivin et al., 2009), SemEval (Kim et al., 2010) and NUS (Nguyen and Kan, 2007). Following previous works, we use the training set of KP20k to train all models. All models are run three times on different seeds. We report the mean and standard deviation of these three runs. Hyperparameters are detailed in Appendix A.1.

5.1 Results on Absent Keyphrase Generation

In Table 1, we show the F1 performance for absent keyphrase generation. As we can see from the ta-

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	R@10	R@50	R@10	R@50	R@10	R@50	R@10	R@50	R@10	R@50
GRU One2Many										
Beam	3.8 ₁	3.8 ₁	5.3 ₂	5.3 ₂	8.2 ₃	8.2 ₃	3.0 ₀	3.0 ₀	8.1 ₁	8.1 ₁
Beam+KPD-R	4.3₃	4.3₄	7.5₄	7.6₃	10.9₁	11.1₂	4.2₃	4.2₃	10.1₁	10.2₁
Beam+KPD-A	4.2 ₂	4.2 ₂	6.4 ₁₀	6.4 ₁₀	9.0 ₃	9.1 ₃	3.9 ₂	4.0 ₂	9.2 ₁	9.3 ₁
GRU One2One										
Beam	6.1 ₃	11.3 ₂	9.5 ₁₀	17.4 ₅	12.3 ₂	22.9 ₁	4.4 ₂	8.1 ₇	14.0 ₀	23.5 ₁
Beam+KPD-R	6.3₃	12.2 ₄	12.1₉	20.3₄	15.6₉	26.8₆	5.6₃	9.4₃	14.6 ₁	24.8 ₃
Beam+KPD-A	5.9 ₄	12.9₇	11.3 ₇	19.0 ₅	14.5 ₇	24.7 ₅	5.0 ₄	9.3 ₃	15.2₁	25.5₁
Transformer One2Set										
Beam	6.7 ₃	12.5 ₃	12.0 ₅	19.6 ₃	13.6 ₁₀	24.4 ₉	5.4 ₂	9.1 ₂	13.5 ₂	23.4 ₂
Beam+KPD-R	5.7 ₁	11.7 ₃	13.3 ₆	21.2 ₅	15.6₅	27.3 ₉	5.9₈	10.4 ₅	14.5 ₀	25.4 ₂
Beam+KPD-A	7.9₃	13.6₄	14.0₂	22.3₈	15.2 ₄	27.7₄	5.9₅	11.0₇	15.6₂	26.6₂

Table 2: Absent keyphrase performance (**Recall**) for different models. KPD represents KPDRop. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1). We bold the best scores per block.

ble, KPDRop always boosts the absent keyphrase performance against a comparable baseline (regardless of whether we use KPDRop-R or KPDRop-A). Comparing greedy with beam search in GRU One2Many, we find that overgenerating with beam search can substantially improve the absent performance over greedy and further applying KPDRop to beam search improves the performance significantly (e.g., on KP20K performance improves from 5.8% to 6.5%). In both GRU One2One and Transformer One2Set, we observe that beam search leads to overgeneration reducing precision and thus the F1@M performance. However, F1@5 generally improves in beam search compared to greedy because only the top 5 keyphrases are considered in this metric. Either way, whether using greedy or beam search, in both One2One models and One2Set models, KPDRop enhances the base performance. Overall, for absent performance, KPDRop-R and KPDRop-A are both competitive against each other.

In Table 2, we show the recall performance for absent keyphrase generation when using high beam size. We find that applying any of the KPDRop techniques substantially improves the recall performance of absent keyphrase generation in all settings and datasets. In GRU One2Many and GRU One2One, KPDRop-R generally performs better than KPDRop-A. Interestingly, in Transformer One2Set KPDRop-A generally performs better than KPDRop-R across all datasets.

The above results validate our intuition that dropping keyphrases in KPDRop forces the models to deeply utilize the context information to learn to predict the dropped keyphrases, and thus, yields

more robust models for absent keyphrase.

5.2 Results on Present Keyphrase Generation

In Table 3, we show the F1 performance for present keyphrase generation. As we hypothesized before, on greedy decoding KPDRop-R can significantly drop the performance of present keyphrases. However, KPDRop-A performs quite competitively against the baselines even in present keyphrase performance. Thus, KPDRop-A can be a balanced approach to boost absent performance without significantly downgrading the present performance. Interestingly, when combined with beam search, even KPDRop-R becomes competitive in present keyphrase generation. While, as we suggested before (in §2.1), KPDRop-R can create a tendency to undergenerate present keyphrases (given that many of them get turned artificially absent), beam search can fight against that tendency through overgeneration. This can sometimes lead to a “sweet spot” when beam search is combined with KPDRop-R making it competitive even in present performance for One2Many and One2One settings.

6 Semi-Supervised Experiments

We also investigate whether KPDRop has something to offer in a low-resource semi-supervised setting. We simulate a semi-supervised setting by randomly splitting the training set of KP20K into two parts. In one part, we keep 5000 samples with their original keyphrases intact, but in the other part, we keep the rest of the data after removing the original keyphrases. Thus, we are left with a low-resource (5000 samples) author-annotated training data, and a huge unlabelled corpus (rest

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5
GRU One2Many										
Greedy	27.4 ₇	27.1 ₈	38.1 ₃	37.7 ₆	34.7 ₉	34.2 ₇	29.6 ₁₀	29.1 ₉	37.3 ₁	37.0 ₀
Greedy+KPD-R	18.7 ₃	18.7 ₃	28.1 ₃	28.1 ₃	28.4 ₆	28.4 ₆	20.1 ₁₄	20.1 ₁₄	30.8 ₂	30.8 ₂
Greedy+KPD-A	25.1 ₄	24.8 ₃	38.5 ₄	38.2 ₃	34.1 ₉	33.5 ₈	27.9 ₇	27.7 ₆	37.0 ₂	36.7 ₂
Beam	38.7 ₃	34.4 ₆	37.0 ₆	41.6 ₅	25.6 ₅	32.8 ₃	31.9 ₁₂	33.4 ₃	32.2 ₁	36.8 ₀
Beam+KPD-R	34.9 ₂	33.2 ₂	42.0 ₄	41.0 ₂	36.7 ₃	36.7 ₂	34.6 ₉	34.1 ₉	36.3 ₁	36.5 ₁
Beam+KPD-A	38.9 ₅	34.6 ₅	39.4 ₅	41.8 ₇	28.4 ₀	34.2 ₈	33.0 ₁₇	33.2 ₇	33.9 ₂	37.1 ₁
GRU One2One										
Beam	25.5 ₁	30.6 ₁	16.5 ₃	40.2 ₁₀	12.4 ₁	29.8 ₁	16.5 ₂	29.2 ₁₆	13.1 ₀	39.4 ₁
Beam+KPD-R	29.6 ₃	30.4 ₃	24.6 ₄	43.9 ₃	17.5 ₁	36.5 ₂	23.6 ₄	32.7 ₁₀	17.0 ₂	37.9 ₀
Beam+KPD-A	26.3 ₄	30.6 ₂	17.1 ₀	38.9 ₆	12.9 ₀	29.4 ₄	17.2 ₅	30.0 ₁₁	13.9 ₁	39.6 ₀
Transformer One2Set										
Greedy	32.2 ₇	31.3 ₆	43.7 ₀	42.1 ₃	35.2 ₅	34.3 ₁₁	34.7 ₂	33.4 ₇	39.2 ₁	37.9 ₄
Greedy+KPD-R	21.1 ₃	21.0 ₃	35.0 ₉	34.9 ₁₂	34.2 ₇	34.2 ₇	26.9 ₇	27.0 ₈	34.9 ₂	34.8 ₂
Greedy+KPD-A	30.6 ₃	29.8 ₄	44.4 ₃	42.6 ₃	35.3 ₆	34.0 ₆	34.4 ₅	33.6 ₅	39.6 ₂	38.5 ₀
Beam	21.7 ₀	32.3 ₄	15.5 ₁	42.3 ₃	11.0 ₁	32.9 ₆	16.3 ₂	33.5 ₈	10.9 ₀	36.4 ₅
Beam+KPD-R	23.2 ₂	27.8 ₈	21.0 ₅	40.0 ₉	14.8 ₈	33.6 ₇	20.9 ₅	32.0 ₅	15.0 ₃	35.1 ₃
Beam+KPD-A	22.7 ₁	32.2 ₆	16.8 ₂	41.8 ₃	11.6 ₂	32.3 ₇	17.6 ₄	34.3 ₁₂	11.7 ₂	36.3 ₀

Table 3: Present keyphrase performance (F1) for different models. KPD represents KPDrop. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1). We bold the best scores per block.

of the KP20K training set). Henceforth, we refer to the former low-resource labelled dataset as LR, and the latter unlabelled corpus as UC.

We investigate a pre-training based approach to utilize UC. Essentially, we pre-train our models first on UC and then fine-tune them on LR. For pre-training, similar to Ye and Wang (2018), we create synthetic labels using an unsupervised keyphrase extraction model. Unlike Ye and Wang (2018), we use a contemporary embedding-based keyphrase extraction model (Liang et al., 2021) to generate the synthetic keyphrases. Particularly, we rank the candidate keyphrases and keep the top 10. Note that in this pre-training setup, the synthetic keyphrases will only be present keyphrases because they are extracted from the input text. This is where KPDROP can make the unsupervised pre-training more interesting by creating artificial absent keyphrases by dropping of the synthetic present keyphrases (and simulating real data). We hypothesize that the application of KPDROP can help our models to learn more effective weights from UC in the self-supervised pre-training stage. We test the effectiveness of using KPDROP to augment self-supervised pre-training by testing the pre-trained model on the labelled test sets after fine-tuning on LR. During fine-tuning on LR, for all models, we always use KPDROP-A as we have already shown this to be beneficial in most supervised contexts. Hyperpa-

rameters are detailed in the Appendix A.1.

6.1 Semi-Supervised Results

For the semi-supervised experiments, report the beam search performance of GRU One2Many for the sake of brevity (greedy performance is in Appendix A.3). In Table 4, we show the absent performance of GRU One2Many in semi-supervised settings³. As we can see from the table, absent performance is near zero in almost all settings. Only when the model is pre-trained (PT) with KPDROP-A or KPDROP-R and then fine-tuned (FT) on LR (PT+KPD-R;FT or PT+KPD-A;FT), there is some degree of absent keyphrase generation. Thus, KPDROP is crucial for downstream absent keyphrase performance in this semi-supervised environment and for domains with low-resource annotated data.

In Table 4, we also show the present performance of GRU One2Many in semi-supervised settings. For present performance, we find that neither training only (FT) on LR nor training only (PT) on UC is as good as pre-training on UC and then fine-tuning the pre-trained model on LR (PT;FT). Although, one exception is the performance on Inspec which can be sometimes better when the model is zero-shot after only training on UC. Either way, we again find that using KPDROP in the pre-training (PT+KPD-R;FT or PT+KPD-A;FT) setting signifi-

³Note that Liang et al. (2021) use an extraction model; thus it has no capabilities for absent keyphrase generation

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	F@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5
Absent Keyphrase Performance										
Unsupervised Extraction Model (Liang et al., 2021)										
	—	0.0	—	0.0	—	0.0	—	0.0	—	0.0
GRU One2Many Models (Beam Search)										
PT	0.0 ₁	0.0 ₁	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₁	0.0 ₁	0.0 ₀	0.0 ₀	0.0 ₀
PT+KPD-R	0.7 ₁	0.7 ₁	0.8 ₂	0.7 ₁	0.6 ₂	0.5 ₂	0.3 ₁	0.2 ₂	0.8 ₁	0.8 ₁
PT+KPD-A	0.1 ₀	0.1 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀
FT	0.7 ₁	0.6 ₁	1.1 ₀	1.1 ₀	0.5 ₀	0.6 ₀	0.3 ₁	0.3 ₁	0.7 ₁	0.7 ₁
PT; FT	0.7 ₂	0.7 ₂	0.5 ₄	0.5 ₄	0.1 ₀	0.1 ₀	0.1 ₁	0.1 ₁	0.4 ₁	0.3 ₁
PT+KPD-A; FT	0.9 ₂	0.9 ₂	2.5 ₃	2.6 ₄	1.8 ₂	1.8 ₂	1.8 ₁	1.8 ₁	1.8 ₁	1.8 ₁
PT+KPD-R; FT	1.7₃	1.7₃	2.9₃	2.8₃	2.8₆	2.8₅	2.0₁	2.0₁	2.7₂	2.7₂
Present Keyphrase Performance										
Unsupervised Extraction Model (Liang et al., 2021)										
	—	32.6	—	20.8	—	18.1	—	13.02	—	17.9
GRU One2Many Models (Beam Search)										
PT	40.9₃	36.0 ₃	20.6 ₂	22.4 ₁₀	18.5 ₂	17.9 ₂	22.8 ₅	22.2 ₁₁	16.0 ₁	17.9 ₂
PT+KPD-R	36.2 ₁₃	34.3 ₄	24.8 ₁₃	23.6 ₂₀	19.2 ₁	18.9 ₄	25.1 ₆	23.2 ₁₁	18.6 ₈	18.0 ₇
PT+KPD-A	40.2 ₃	35.7 ₄	21.0 ₆	22.8 ₅	18.8 ₅	19.8 ₂	23.4 ₆	22.7 ₆	16.1 ₂	17.9 ₂
FT	27.8 ₁₁	26.5 ₉	32.6 ₂	31.8 ₃	26.6 ₈	26.4 ₈	27.8 ₈	26.7 ₇	27.1 ₄	26.9 ₄
PT; FT	36.8 ₁₇	33.5 ₂₅	33.8 ₃	33.5 ₂	26.2 ₁₂	27.5 ₆	29.1 ₅	27.6 ₇	26.2 ₃	27.2 ₂
PT+KPD-A; FT	39.9 ₅	36.0 ₁₀	36.3₄	35.7₂	29.2 ₄	29.6 ₃	32.0₄	30.8₁₂	28.2 ₃	28.8 ₂
PT+KPD-R; FT	40.0 ₁₁	36.9₁₀	36.2 ₁₁	34.6 ₉	31.5₃	31.1₇	30.9 ₅	29.9 ₁₄	29.7₇	29.8₈

Table 4: Absent and present keyphrase performance using Beam Search for GRU One2Many models in a semi-supervised setup. KPD represents KPDrop. PT represents pre-training on the synthetic data (UC). PT+KPD-R or PT+KPD-A represents pre-training on UC with KPD-R or KPD-A respectively. FT represents fine tuning or training on the low resource labelled data (LR). PT (or PT+KPD-A or PT+KPD-R) followed by “; FT” represents that the pre-training was followed by fine-tuning of the pre-trained model on LR. We bold the overall best scores. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1).

cantly boosts present performance after fine-tuning.

Thus, using KPDROP to make the pre-training stage more challenging by pushing the model to predict missing present keyphrases helps not only with absent keyphrase performance but also with present keyphrase performance after fine-tuning. In between KPDROP-R and KPDROP-A, the former generally performs better in the pre-training stage. Thus, during pre-training it is better to *replace* synthetic labels of fully present keyphrases with its KPDropped version. In Appendix A.3, we also observe similar patterns from other models (One2One and One2Set) in semi-supervised settings.

7 Preliminary Experiments on Pre-trained Models

We also did a few preliminary experiments on applying KPDROP-R to a large pre-trained Seq2Seq language model, in particular, T5 (Raffel et al., 2020). We present the results in Appendix A.4. Consistent with previous results, we find that KPDROP-R increases absent performance for T5 as well. However, overall, we found the perfor-

mance of T5 baseline to be limited compared to trained-from-scratch models like Transformer One2Set. Similarly, other reported performances on pre-trained models have been generally lower than Transformers One2Set too. For example, even after large scale specialized pre-training for keyphrases, Kulkarni et al. (2022) reports only comparable performance on present keyphrases to Transformer One2Set (Ye et al., 2021b) and much less in absent performance. On the other hand, using pre-trained models, Wu et al. (2021)⁴ achieve comparable on absent performance with One2Set under greedy search, but much less in present performance. However, it should not be too difficult to adapt a pre-trained model into a one2set framework during fine-tuning. This can be a promising future direction and form a stronger base model for KPDROP.

8 Related Work

There is a wide variety of approaches (Hasan and Ng, 2014; Caragea et al., 2014; Das Gollapalli and

⁴We are referring to the latest ArXiv version (v2) which holds the globally latest version of the paper.

Caragea, 2014; Gollapalli and Li, 2016; Sterckx et al., 2016; Florescu and Caragea, 2017; Zhang et al., 2017; Boudin, 2018; Mahata et al., 2018; Sun et al., 2019; Al-Zaidy et al., 2019; Campos et al., 2020; Santosh et al., 2020; Sahrawat et al., 2020; Sun et al., 2020; Song et al., 2021; Patel and Caragea, 2021) for keyphrase extraction exclusively. Meng et al. (2017) diverged from pure extractive methods by introducing a seq2seq model (CopyRNN) for generation of both present and absent keyphrases. Chen et al. (2018) extended CopyRNN with keyphrase correlation constraints and Zhao and Zhang (2019) extended it with linguistic constraints. Ye and Wang (2018); Wu et al. (2022) investigated keyphrase generation (KG) in semi-supervised or resource-constrained settings. Chen et al. (2019b) used a title-guided encoding method for better KG. Wang et al. (2019) incorporated a topic-model to enhance KG. Yuan et al. (2020) extended CopyRNN by introducing the CatSeq model that can generate a concatenation of dynamically determined variable number of keyphrases. Chan et al. (2019); Luo et al. (2021) improved KG using reinforcement learning whereas Swaminathan et al. (2020a,b); Lancioni et al. (2020) do so using GANs. A few approaches (Chen et al., 2019a; Diao et al., 2020; Kim et al., 2021; Ye et al., 2021a) augmented KG with information from retrieved documents.

Multiple approaches (Chen et al., 2019a; Liu et al., 2020; Ahmad et al., 2021; Zhao et al., 2021; Wu et al., 2021) took a joint-training or multi-tasking approach to do both present keyphrase extraction and absent keyphrase generation. Chen et al. (2020) and Ye et al. (2021b) changed the decoder to better respect the structure of keyphrases. Luo et al. (2020) changed the encoder to better respect the input document structure. Huang et al. (2021) proposed a new beam-search-based adaptive decoding method. Meng et al. (2021); Kulkarni et al. (2022) investigated pre-training objectives for KG. Both, however, rely on labelled pre-training data.

9 Conclusion

Our proposal, KPDRÖP, randomly drops present keyphrases from a document to turn them artificially absent. This encourages the model to learn to better exploit the context in the input to be able to infer keyphrases that are absent from the text but otherwise topically relevant. The results show that KPDRÖP serves as a simple model-agnostic

method to substantially improve absent (and sometimes, present) keyphrase performance in both supervised and semi-supervised (low resource) settings when large annotated datasets for keyphrase generation are not available. In future, we would like to explore integration of KPDRÖP with large-scale pre-training.

10 Limitations

KPDRÖP is a simple yet effective approach to improving performance of keyphrase generation in both large datasets and low resource datasets, which makes it applicable to a wide range of domains where keyphrases are necessary. However, one limitation of KPDRÖP (especially KPDRÖP-A) is that it can increase the computation cost during training because both the effective mini-batch size and the effective training dataset size per epoch is doubled through data augmentation. Yet, most data augmentation techniques share the same limitation.

In addition, KPDRÖP-R can potentially harm the performance of present keyphrases in some contexts (especially when using greedy search in a supervised setting). To address this, we can simply use the absent predictions of the model trained with KPDRÖP and the present predictions of the model trained without KPDRÖP.

11 Ethics Statement

Our technique is specifically designed to improve keyphrase generation. Keyphrase generation is a well-established traditional NLP task that is useful in several application contexts related to organization of information. We do not foresee any immediate ethical concern following from our contribution in this area.

Acknowledgements

This research is supported in part by NSF CAREER award #1802358, NSF CRI award #1823292, NSF IIS award #2107518, and UIC Discovery Partners Institute (DPI) award. Any opinions, findings, and conclusions expressed here are those of the authors and do not necessarily reflect the views of NSF or DPI. We thank AWS for computational resources used for this study. We also thank our anonymous reviewers for their constructive feedback and suggestions.

References

- Wasi Ahmad, Xiao Bai, Soomin Lee, and Kai-Wei Chang. 2021. Select, extract and generate: Neural keyphrase generation with layer-wise coverage attention. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1389–1404.
- Rabah Al-Zaidy, Cornelia Caragea, and C. Lee Giles. 2019. Bi-lstm-crf sequence labeling for keyphrase extraction from scholarly documents. In *Proceedings of The Web Conference (WWW 2019), San Francisco, California, USA*.
- Rohan Anil, Vineet Gupta, Tomer Koren, and Yoram Singer. 2019. Memory efficient adaptive optimization. *Advances in Neural Information Processing Systems*, 32:9749–9758.
- Haoli Bai, Zhuangbin Chen, Michael R. Lyu, Irwin King, and Zenglin Xu. 2018. [Neural relational topic models for scientific article analysis](#). In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM '18*, page 27–36, New York, NY, USA. Association for Computing Machinery.
- Florian Boudin. 2018. [Unsupervised keyphrase extraction with multipartite graphs](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 667–672, New Orleans, Louisiana. Association for Computational Linguistics.
- Florian Boudin and Ygor Gallina. 2021. [Redefining absent keyphrases and their effect on retrieval effectiveness](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4185–4193, Online. Association for Computational Linguistics.
- Florian Boudin, Ygor Gallina, and Akiko Aizawa. 2020. [Keyphrase generation for scientific document retrieval](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1118–1126, Online. Association for Computational Linguistics.
- Ricardo Campos, Vítor Mangaravite, Arian Pasquali, Alípio Jorge, Célia Nunes, and Adam Jatowt. 2020. [Yake! keyword extraction from single documents using multiple local features](#). *Information Sciences*, 509:257–289.
- Cornelia Caragea, Florin Adrian Bulgarov, Andreea Godea, and Sujatha Das Gollapalli. 2014. [Citation-enhanced keyphrase extraction from research papers: A supervised approach](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1435–1446, Doha, Qatar. Association for Computational Linguistics.
- Hou Pong Chan, Wang Chen, Lu Wang, and Irwin King. 2019. [Neural keyphrase generation via reinforcement learning with adaptive rewards](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2163–2174, Florence, Italy. Association for Computational Linguistics.
- Jun Chen, Xiaoming Zhang, Yu Wu, Zhao Yan, and Zhoujun Li. 2018. [Keyphrase generation with correlation constraints](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4057–4066, Brussels, Belgium. Association for Computational Linguistics.
- Wang Chen, Hou Pong Chan, Piji Li, Lidong Bing, and Irwin King. 2019a. [An integrated approach for keyphrase generation via exploring the power of retrieval and extraction](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2846–2856, Minneapolis, Minnesota. Association for Computational Linguistics.
- Wang Chen, Hou Pong Chan, Piji Li, and Irwin King. 2020. [Exclusive hierarchical decoding for deep keyphrase generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1095–1105, Online. Association for Computational Linguistics.
- Wang Chen, Yifan Gao, Jiani Zhang, Irwin King, and Michael R Lyu. 2019b. [Title-guided encoding for keyphrase generation](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6268–6275.
- Sujatha Das Gollapalli and Cornelia Caragea. 2014. [Extracting keyphrases from research papers using citation networks](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 28(1).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Shizhe Diao, Yan Song, and Tong Zhang. 2020. [Keyphrase generation with cross-document attention](#). *ArXiv*.
- Angela Fan, Edouard Grave, and Armand Joulin. 2020. [Reducing transformer depth on demand with structured dropout](#). In *International Conference on Learning Representations*.
- Corina Florescu and Cornelia Caragea. 2017. [PositionRank: An unsupervised approach to keyphrase extraction from scholarly documents](#). In *Proceedings*

- of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1105–1115, Vancouver, Canada. Association for Computational Linguistics.
- George W. Furnas, Thomas K. Landauer, Louis M. Gomez, and Susan T. Dumais. 1987. The vocabulary problem in human-system communication. *Communications of the ACM*, 30(11):964–971.
- Sujatha Das Gollapalli and Xiaoli Li. 2016. [Keyphrase extraction using sequential labeling](#). *ArXiv*, abs/1608.00329.
- Kazi Saidul Hasan and Vincent Ng. 2014. [Automatic keyphrase extraction: A survey of the state of the art](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1262–1273, Baltimore, Maryland. Association for Computational Linguistics.
- Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q. Weinberger. 2016. Deep networks with stochastic depth. In *Computer Vision – ECCV 2016*, pages 646–661, Cham. Springer International Publishing.
- Xiaoli Huang, Tongge Xu, Lvan Jiao, Yueran Zu, and Youmin Zhang. 2021. [Adaptive beam search decoding for discrete keyphrase generation](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(14):13082–13089.
- Anette Hulth. 2003. [Improved automatic keyword extraction given more linguistic knowledge](#). In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*, pages 216–223.
- Mohit Iyyer, Varun Manjunatha, Jordan Boyd-Graber, and Hal Daumé III. 2015. [Deep unordered composition rivals syntactic methods for text classification](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1681–1691, Beijing, China. Association for Computational Linguistics.
- James M. Kates and Kathryn H. Arehart. 2014. [The hearing-aid speech perception index \(haspi\)](#). *Speech Communication*, 65:75–93.
- Jihyuk Kim, Myeongho Jeong, Seungtaek Choi, and Seung-won Hwang. 2021. [Structure-augmented keyphrase generation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2657–2667, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Su Nam Kim, Olena Medelyan, Min-Yen Kan, and Timothy Baldwin. 2010. [SemEval-2010 task 5 : Automatic keyphrase extraction from scientific articles](#). In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 21–26, Uppsala, Sweden. Association for Computational Linguistics.
- Mikalai Krapivin, Aliaksandr Autaeu, and Maurizio Marchese. 2009. Large dataset for keyphrases extraction. *University of Trento*.
- H. W. Kuhn. 1955. [The hungarian method for the assignment problem](#). *Naval Research Logistics Quarterly*, 2(1-2):83–97.
- Mayank Kulkarni, Debanjan Mahata, Ravneet Arora, and Rajarshi Bhowmik. 2022. [Learning rich representation of keyphrases from text](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 891–906, Seattle, United States. Association for Computational Linguistics.
- Giuseppe Lancioni, Saida S.Mohamed, Beatrice Portelli, Giuseppe Serra, and Carlo Tasso. 2020. [Keyphrase generation with GANs in low-resources scenarios](#). In *Proceedings of SustaiNLP: Workshop on Simple and Efficient Natural Language Processing*, pages 89–96, Online. Association for Computational Linguistics.
- Xinnian Liang, Shuangzhi Wu, Mu Li, and Zhoujun Li. 2021. [Unsupervised keyphrase extraction by jointly modeling local and global context](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 155–164, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Rui Liu, Zheng Lin, and Weiping Wang. 2020. [Keyphrase prediction with pre-trained language model](#). *ArXiv*, abs/2004.10462.
- Yichao Luo, Zhengyan Li, Bingning Wang, Xiaoyu Xing, Qi Zhang, and Xuanjing Huang. 2020. [Sensenet: Neural keyphrase generation with document structure](#). *ArXiv*, abs/2012.06754.
- Yichao Luo, Yige Xu, Jiacheng Ye, Xipeng Qiu, and Qi Zhang. 2021. [Keyphrase generation with fine-grained evaluation-guided reinforcement learning](#). In *Proceedings of EMNLP findings*.
- Debanjan Mahata, John Kuriakose, Rajiv Ratn Shah, and Roger Zimmermann. 2018. [Key2Vec: Automatic ranked keyphrase extraction from scientific articles using phrase embeddings](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 634–639, New Orleans, Louisiana. Association for Computational Linguistics.
- Rui Meng, Xingdi Yuan, Tong Wang, Sanqiang Zhao, Adam Trischler, and Daqing He. 2021. [An empirical study on neural keyphrase generation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4985–5007, Online. Association for Computational Linguistics.
- Rui Meng, Sanqiang Zhao, Shuguang Han, Daqing He, Peter Brusilovsky, and Yu Chi. 2017. [Deep keyphrase](#)

- generation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 582–592.
- Xinfan Meng, Furu Wei, Xiaohua Liu, Ming Zhou, Sujian Li, and Houfeng Wang. 2012. Entity-centric topic-oriented opinion summarization in twitter. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 379–387.
- Thuy Dung Nguyen and Min-Yen Kan. 2007. Keyphrase extraction in scientific publications. In *Asian Digital Libraries. Looking Back 10 Years and Forging New Frontiers*, pages 317–326, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Krutarth Patel and Cornelia Caragea. 2021. [Exploiting position and contextual word embeddings for keyphrase extraction from scientific papers](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1585–1591, Online. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67.
- Dhruva Sahrawat, Debanjan Mahata, Haimin Zhang, Mayank Kulkarni, Agniv Sharma, Rakesh Gosangi, Amanda Stent, Yaman Kumar Singla, Rajiv Ratn Shah, and Roger Zimmermann. 2020. Keyphrase extraction as sequence labeling using contextualized embeddings. *Information Retrieval, 42nd European Conference on IR Research, ECIR 2020*, 12036:328–335.
- T.y.s.s Santosh, Debarshi Kumar Sanyal, Plaban Kumar Bhowmick, and Partha Pratim Das. 2020. [SaSAKE: Syntax and semantics aware keyphrase extraction from research papers](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5372–5383, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Mingyang Song, Liping Jing, and Lin Xiao. 2021. [Importance Estimation from Multiple Perspectives for Keyphrase Extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2726–2736, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. [Dropout: A simple way to prevent neural networks from overfitting](#). *Journal of Machine Learning Research*, 15(56):1929–1958.
- Lucas Sterckx, Cornelia Caragea, Thomas Demeester, and Chris Develder. 2016. [Supervised keyphrase extraction as positive unlabeled learning](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1924–1929, Austin, Texas. Association for Computational Linguistics.
- Si Sun, Zhenghao Liu, Chenyan Xiong, Zhiyuan Liu, and Jie Bao. 2020. [Joint keyphrase chunking and salience ranking with bert](#). *ArXiv*.
- Zhiqing Sun, Jian Tang, Pan Du, Zhi-Hong Deng, and Jian-Yun Nie. 2019. [Divgraphpointer: A graph pointer network for extracting diverse keyphrases](#). In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR’19*, page 755–764, New York, NY, USA. Association for Computing Machinery.
- Avinash Swaminathan, Raj Kuwar Gupta, Haimin Zhang, Debanjan Mahata, Rakesh Gosangi, and Rajiv Ratn Shah. 2020a. [Keyphrase generation for scientific articles using gans \(student abstract\)](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(10):13931–13932.
- Avinash Swaminathan, Haimin Zhang, Debanjan Mahata, Rakesh Gosangi, Rajiv Ratn Shah, and Amanda Stent. 2020b. [A preliminary exploration of GANs for keyphrase generation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8021–8030, Online. Association for Computational Linguistics.
- Yue Wang, Jing Li, Hou Pong Chan, Irwin King, Michael R. Lyu, and Shuming Shi. 2019. [Topic-aware neural keyphrase generation for social media language](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2516–2526, Florence, Italy. Association for Computational Linguistics.
- Thomas Wolf, Julien Chaumond, Lysandre Debut, Victor Sanh, Clement Delangue, Anthony Moi, Pierric Cistac, Morgan Funtowicz, Joe Davison, Sam Shleifer, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.
- Di Wu, Wasi Uddin Ahmad, Sunipa Dev, and Kai-Wei Chang. 2022. [Representation learning for resource-constrained keyphrase generation](#). *ArXiv*, abs/2203.08118.
- Huanqin Wu, Wei Liu, Lei Li, Dan Nie, Tao Chen, Feng Zhang, and Di Wang. 2021. [UniKeyphrase: A unified extraction and generation framework for keyphrase prediction](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 825–835, Online. Association for Computational Linguistics.

- Hai Ye and Lu Wang. 2018. [Semi-supervised learning for neural keyphrase generation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4142–4153, Brussels, Belgium. Association for Computational Linguistics.
- Jiacheng Ye, Ruijian Cai, Tao Gui, and Qi Zhang. 2021a. [Heterogeneous graph neural networks for keyphrase generation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2705–2715, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jiacheng Ye, Tao Gui, Yichao Luo, Yige Xu, and Qi Zhang. 2021b. [One2Set: Generating diverse keyphrases as a set](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4598–4608, Online. Association for Computational Linguistics.
- Xingdi Yuan, Tong Wang, Rui Meng, Khushboo Thaker, Peter Brusilovsky, Daqing He, and Adam Trischler. 2020. One size does not fit all: Generating and evaluating variable number of keyphrases. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7961–7975.
- Yuxiang Zhang, Yaocheng Chang, Xiaoqing Liu, Sujatha Das Gollapalli, Xiaoli Li, and Chunjing Xiao. 2017. [Mike: Keyphrase extraction by integrating multidimensional information](#). In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM '17*, page 1349–1358, New York, NY, USA. Association for Computing Machinery.
- Jing Zhao, Junwei Bao, Yifan Wang, Youzheng Wu, Xiaodong He, and Bowen Zhou. 2021. [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*, pages 5717–5726, Online. Association for Computational Linguistics.
- Jing Zhao and Yuxiang Zhang. 2019. [Incorporating linguistic constraints into keyphrase generation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5224–5233, Florence, Italy. Association for Computational Linguistics.
- Özlem Ünlü and Aydın Çetin. 2019. [A survey on keyword and key phrase extraction with deep learning](#). In *2019 3rd International Symposium on Multidisciplinary Studies and Innovative Technologies (ISM-SIT)*, pages 1–6.

A Appendix

A.1 Hyperparameters

For GRU One2Many (CatSeq) we use hyperparameters similar to prior works (Chan et al., 2019). Like prior works (Yuan et al., 2020; Meng et al., 2021), when using beam search on one2many models we use a beam size of 50. We use the same model hyperparameters of GRU One2One models as used for GRU One2Many. However, during beam search we use a beam size of 200 as is standard for one2one models in prior work (Meng et al., 2017, 2021). A lower beam size (50 instead of 200) is used for one2many models because they can generate multiple keyphrases per beam. So a lower size is used to make them more comparable with one2one models with higher beam size. For Transformer One2Set model we use the same hyperparameters as Ye et al. (2021b). Given the similarity of One2Set models with One2Many models in their ability to generate multiple keyphrases per beam, we also use a beam size of 50 for Transformer One2Set. Anytime we use beam search, we also use length normalization on beam search with a length co-efficient of 0.8. We use the same settings during semi-supervised pre-training or fine-tuning. We tested KPDRP-R rates among {0.3, 0.5, 0.7, 0.9, 1.0} on GRU One2Many over the validation set after one epoch training on the full KP20K training set. We found 0.7 to be the best performing rate for validation absent performance. We use this same rate for KPDRP-A and other models in both supervised and semi-supervised settings. All the experiments were done in a single Nvidia RTX A6000.

A.2 More Evaluations

Chan et al. (2019) modified the original $F_1@5$ (as used in the main paper and other prior works (Meng et al., 2017; Yuan et al., 2020)) such that the denominator in the precision is always set to 5 even when the total predictions are less than 5. To avoid confusion and better distinguish from the original $F_1@5$ we refer to the modified metric as $F_1@5C$. Throughout the paper, for the sake of brevity we mainly report performance in $F_1@5$ instead of $F_1@5C$. This is because although $F_1@5C$ achieve the goal of differentiating itself from $F_1@M$ reports in greedy search contexts, it can be a little misleading. For example, $F_1@5C$ can artificially penalize the model for predicting less than 5 keyphrases even when the ground truth itself contains less than

5 keyphrases. Nevertheless, in Table 5, we still report the $F_1@5C$ present and absent performance of our models from our supervised experiments for the sake of better comparison with prior works that use $F_1@5C$.

A.3 More Semi-Supervised Experiments

In Table 6 we present the greedy decoding performance of GRU One2Many models in semi-supervised settings. In Table 7 we present the performance of GRU One2One models (beam search) in semi-supervised settings. In Table 8 and 9 we present the greedy decoding performance and beam decoding performance of Transformer One2Set models in semi-supervised settings respectively. Generally, we notice similar patterns across all the models as were found and discussed for GRU One2Many models in §6.1. Overall, KPDRP-R consistently serves as a crucial ingredient in the pre-training stage to enable substantially improved downstream performance in both present and absent keyphrase generation.

A.4 Preliminary Experiments on Additional Baselines

We also present the results of applying KPDRP-R on a large pre-trained model T5 (Raffel et al., 2020), and another specialized KG model, ExHiRD (Chen et al., 2020) in Table 10. In both cases, we see the promise of KPDRP-R in improving the absent performance.

Below we describe the hyperparameters that were used for the preliminary experiments described in this section.

For ExHiRD, we use the same hyperparameters as in the original paper (Chen et al., 2020)⁵. For T5, we use a maximum of 10 epochs, early stopping with a patience of 2 (the training is terminated if validation loss does not improve for 2 consecutive epochs), a batch size as 64, a maximum gradient norm clipping of 5, and SM3 (Anil et al., 2019) as the optimizer. The initial learning rate for T5 was set to be 0.1. We apply learning rate (lr) warmup as follows⁶:

$$lr_s = lr_0 \cdot \min(1, (s/w)^2) \quad (1)$$

lr_s indicates the learning rate at step s . lr_0 is the initial learning rate (0.1). s indicates the current

⁵<https://github.com/Chen-Wang-CUHK/ExHiRD-DKG>

⁶based on the recommended procedure for SM3 (<https://github.com/google-research/google-research/tree/master/sm3>).

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	F1@5C		F1@5C		F1@5C		F1@5C		F1@5C	
	Pre	Abs	Pre	Abs	Pre	Abs	Pre	Abs	Pre	Abs
GRU One2Many										
Greedy	22.1 ₆	0.7 ₁	31.2 ₃	1.6 ₃	26.4 ₅	1.8 ₃	24.6 ₈	1.6 ₂	29.3 ₀	1.6 ₀
Greedy+KPD-R	12.5 ₃	0.6 ₀	18.4 ₁	1.7 ₃	15.6 ₃	2.6 ₁	13.9 ₁₂	1.8 ₁	17.7 ₁	1.9 ₀
Greedy+KPD-A	19.7 ₄	0.8 ₁	30.6 ₆	1.8 ₃	24.7 ₂	2.4 ₁	23.2 ₆	1.8 ₃	28.1 ₅	2.0 ₂
Beam	34.1 ₅	2.3 ₀	41.3 ₆	4.0 ₂	32.6 ₃	5.2 ₂	33.3 ₃	3.2 ₁	36.6 ₀	4.4 ₁
Beam+KPD-R	30.3 ₂	2.6 ₂	37.8 ₃	5.7 ₁	31.9 ₂	6.6 ₂	32.4 ₁₃	4.3 ₄	32.3 ₁	5.4 ₀
Beam+KPD-A	34.0 ₅	2.5 ₁	41.6 ₇	5.0 ₆	33.8 ₉	5.5 ₂	33.1 ₇	4.3 ₃	36.6 ₀	5.0 ₀
GRU One2One										
Beam	30.5 ₂	2.8 ₁	40.2 ₁₀	5.7 ₇	29.8 ₁	5.9 ₃	29.2 ₁₆	3.8 ₁	39.4 ₁	6.2 ₀
Beam+KPD-R	30.4 ₃	2.9 ₃	43.9 ₃	7.5 ₃	36.5 ₂	7.8 ₂	32.7 ₁₀	4.9 ₅	37.9 ₀	6.5 ₁
Beam+KPD-A	30.6 ₂	2.7 ₁	38.9 ₆	6.7 ₂	29.4 ₄	6.5 ₃	30.0 ₁₁	4.1 ₃	39.6 ₀	6.8 ₀
Transformer One2Set										
Greedy	27.6 ₅	1.9 ₃	38.8 ₄	4.3 ₅	30.9 ₉	4.1 ₂	31.2 ₁	2.6 ₂	34.4 ₃	3.3 ₀
Greedy+KPD-R	15.3 ₂	2.0 ₀	26.0 ₉	5.3 ₅	22.2 ₄	5.9 ₅	20.7 ₇	3.9 ₃	24.0 ₂	4.5 ₁
Greedy+KPD-A	25.7 ₃	2.1 ₁	38.0 ₂	5.2 ₈	29.5 ₇	4.6 ₄	30.3 ₇	3.6 ₁	33.9 ₃	4.2 ₁
Beam	32.3 ₄	3.3 ₂	42.3 ₃	7.0 ₅	32.9 ₆	6.7 ₅	33.5 ₈	4.7 ₄	36.4 ₅	5.8 ₀
Beam+KPD-R	27.8 ₈	2.4 ₁	40.0 ₉	7.2 ₅	33.7 ₇	7.3 ₂	32.0 ₅	5.2 ₄	35.1 ₃	6.1 ₀
Beam+KPD-A	32.2 ₆	3.6 ₂	41.8 ₃	7.9 ₄	32.3 ₇	7.8 ₂	34.3 ₁₂	5.3 ₄	36.3 ₀	6.7 ₀

Table 5: Present and absent keyphrase performance for different models. Pre represents present performance and Abs represents absent performance. KPD represents KPDrop. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1). We bold the best scores per block.

update step number. w indicates total warmup steps (set as 2000). The initial learning rate (0.1) was tuned using grid search based on validation loss among the following choices: {0.1, 0.01, 0.001}. For each trial during hyperparameter optimization we use a maximum of 1 epoch. We only tune the baselines (where KPDROP is unapplied). We do not separately tune other hyperparameters when KPDROP is applied.

We use the T5 implementation as provided by (Wolf et al., 2020). Both ExHiRD and T5 models are trained using teacher forcing mechanism. During inference, we set the maximum phrase length as 50 for both. Keyphrase tokens are greedily generated during inference. A KPDROP rate of 0.7 was used - same as the other previous experiments.

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	F@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5
GRU One2Many Models (Greedy Search)										
Absent Keyphrase Performance										
PT	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀
PT+KPD-R	0.2 ₁	0.2 ₁	0.2 ₂	0.2 ₂	0.1 ₀	0.1 ₀	0.0 ₀	0.0 ₀	0.3 ₁	0.3 ₁
PT+KPD-A	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀
FT	0.1 ₁	0.1 ₁	0.0 ₀	0.1 ₀	0.2 ₀	0.2 ₀	0.3 ₀	0.3 ₀	0.1 ₀	0.1 ₀
PT; FT	0.0 ₁	0.0 ₁	0.2 ₁	0.2 ₁	0.0 ₀	0.0 ₀	0.1 ₁	0.1 ₁	0.1 ₀	0.1 ₀
PT+KPD-A; FT	0.2 ₁	0.2 ₁	0.8₄	0.8₄	0.4 ₁	0.4 ₁	0.7₀	0.7₀	0.6 ₀	0.6 ₀
PT+KPD-R; FT	0.8₃	0.8₃	0.5 ₂	0.5 ₂	0.5₀	0.5₀	0.4 ₂	0.4 ₂	0.8₁	0.8₁
Present Keyphrase Performance										
PT	35.0 ₇	34.2 ₆	23.7 ₂	23.7 ₂	20.1 ₄	20.3 ₂	24.5 ₆	24.8₅	18.4 ₁	18.7 ₁
PT+KPD-R	32.0 ₄	32.0 ₅	25.3 ₃	24.8 ₅	20.5 ₅	20.3 ₃	25.0₁₁	24.7 ₁₁	19.3 ₂	19.0 ₃
PT+KPD-A	35.5₅	34.4₄	23.8 ₉	23.9 ₇	19.3 ₆	20.0 ₃	24.5 ₁₅	24.3 ₁₀	18.2 ₄	18.5 ₃
FT	14.3 ₁₂	14.3 ₁₂	24.7 ₂₁	24.7 ₂₁	23.2 ₁₂	23.2 ₁₂	15.8 ₂	15.8 ₂	22.8 ₈	22.8 ₈
PT; FT	20.9 ₁₈	20.9 ₁₈	27.2 ₁₆	27.2 ₁₆	25.5 ₁₀	25.5 ₁₀	21.0 ₅	21.0 ₅	25.1 ₇	25.1 ₇
PT+KPD-A; FT	22.0 ₃	22.0 ₃	30.1₇	30.1₇	26.5 ₇	26.5 ₇	23.0 ₆	23.0 ₆	27.0 ₂	27.0 ₂
PT+KPD-R; FT	20.5 ₄	20.5 ₄	30.1₈	30.1₈	27.8₄	27.8₄	22.3 ₁₂	22.3 ₁₂	27.1₂	27.1₂

Table 6: Absent and present keyphrase performance using Greedy Search for GRU One2Many models in a semi-supervised setup. KPD represents KPDrop. PT represents pre-training on the synthetic data (UC). PT+KPD-R or PT+KPD-A represents pre-training on UC with KPD-R or KPD-A respectively. FT represents fine tuning or training on the low resource labelled data (LR). PT (or PT+KPD-A or PT+KPD-R) followed by “; FT” represents that the pre-training was followed by fine-tuning of the pre-trained model on LR. We bold the overall best scores. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1).

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	F@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5
GRU One2One (Beam Search)										
Absent Keyphrase Performance										
PT	0.0 ₀	0.2 ₁	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.0 ₀	0.1 ₁	0.0 ₀	0.1 ₀
PT+KPD-R	0.4 ₀	1.2 ₁	0.7 ₀	0.9 ₂	0.6 ₀	1.6 ₂	0.9 ₁	0.9 ₂	0.4 ₀	1.3 ₂
PT+KPD-A	0.1 ₀	0.2 ₁	0.1 ₀	0.1 ₁	0.0 ₀	0.0 ₀	0.1 ₀	0.0 ₀	0.1 ₀	0.1 ₀
FT	0.2 ₀	0.3 ₀	0.6 ₀	1.2 ₂	0.3 ₀	0.9 ₁	0.4 ₀	0.4 ₀	0.3 ₀	0.8 ₁
PT; FT	0.2 ₀	0.6 ₁	0.6 ₁	1.3 ₃	0.4 ₀	1.0 ₁	0.3 ₁	0.5 ₃	0.3 ₀	0.8 ₁
PT+KPD-A; FT	0.3 ₀	1.0 ₁	0.7 ₀	1.4 ₂	0.5 ₀	1.4 ₁	0.6 ₁	1.1 ₃	0.4 ₀	1.3 ₀
PT+KPD-R; FT	0.6₀	1.9₂	1.4₀	4.7₂	1.1₀	4.4₂	1.2₀	3.0₂	0.8₀	3.7₀
Present Keyphrase Performance										
PT	14.9 ₁	34.1 ₈	9.3 ₀	25.1 ₄	6.7 ₀	20.3 ₆	11.0 ₁	25.2 ₃	6.6 ₁	19.0 ₁
PT+KPD-R	21.6₆	33.6 ₅	12.6 ₃	25.3 ₆	9.8 ₁	19.8 ₂	14.6 ₇	23.1 ₈	9.2 ₁	18.4 ₂
PT+KPD-A	15.7 ₂	34.5₅	9.7 ₁	25.9 ₇	7.0 ₁	20.3 ₅	11.1 ₂	25.4 ₁₁	6.8 ₁	19.0 ₅
FT	20.6 ₁₃	20.5 ₂₂	17.8₁₅	31.9 ₂₀	12.5₉	23.4 ₁₇	16.2₁₈	21.8 ₂₄	13.5₈	25.2 ₈
PT; FT	20.6 ₃	26.1 ₇	13.7 ₅	36.6 ₅	10.1 ₁	28.2 ₆	13.9 ₃	28.1 ₁₅	10.0 ₂	29.5 ₁
PT+KPD-A; FT	20.6 ₁₀	27.8 ₁₀	13.8 ₈	35.2 ₁₇	10.1 ₆	28.9 ₉	14.0 ₇	28.7 ₈	9.9 ₆	29.2 ₄
PT+KPD-R; FT	21.4 ₃	29.4 ₈	14.3 ₂	39.0₈	10.5 ₂	30.7₄	15.0 ₁	30.4₄	10.4 ₂	32.5₁

Table 7: Absent and present keyphrase performance using Beam Search for GRU One2One models in a semi-supervised setup. KPD represents KPDrop. PT represents pre-training on the synthetic data (UC). PT+KPD-R or PT+KPD-A represents pre-training on UC with KPD-R or KPD-A respectively. FT represents fine tuning or training on the low resource labelled data (LR). PT (or PT+KPD-A or PT+KPD-R) followed by “; FT” represents that the pre-training was followed by fine-tuning of the pre-trained model on LR. We bold the overall best scores. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1).

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	F@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5
Transformer One2Set (Greedy Search)										
Absent Keyphrase Performance										
PT	0.2 ₁	0.2 ₁	0.1 ₁	0.1 ₁	0.3 ₁	0.3 ₁	0.0 ₀	0.0 ₀	0.1 ₀	0.1 ₀
PT+KPD-R	0.6 ₀	0.6 ₀	0.8 ₁	0.8 ₁	0.8 ₅	0.7 ₅	0.3 ₂	0.3 ₂	0.3 ₀	0.6 ₀
PT+KPD-A	0.3 ₁	0.3 ₁	0.2 ₁	0.2 ₁	0.2 ₀	0.2 ₀	0.1 ₁	0.1 ₁	0.2 ₀	0.2 ₀
FT	0.2 ₃	0.2 ₃	0.4 ₃	0.4 ₃	0.7 ₃	0.7 ₃	0.3 ₁	0.3 ₁	0.5 ₁	0.5 ₁
PT; FT	0.3 ₂	0.3 ₂	0.5 ₁	0.5 ₁	0.4 ₂	0.4 ₂	0.0 ₀	0.0 ₀	0.4 ₂	0.4 ₂
PT+KPD-A; FT	0.5 ₂	0.5 ₂	0.7 ₂	0.7 ₂	1.1 ₁	1.1 ₁	1.1 ₅	1.1 ₅	1.2 ₀	1.2 ₀
PT+KPD-R; FT	1.1₁	1.1₁	2.9₉	2.9₉	2.5₅	2.5₅	2.0₆	2.0₆	2.4₁	2.4₁
Present Keyphrase Performance										
PT	34.6 ₅	28.6 ₁₇	25.9 ₃	23.3 ₁₀	20.1 ₄	17.8 ₄	27.1 ₄	23.3 ₇	18.9 ₁	17.2 ₄
PT+KPD-R	35.7₆	31.7₁₃	26.3 ₉	23.9 ₁₁	20.4 ₁₀	18.3 ₅	26.0 ₇	23.0 ₁₈	19.3 ₈	17.8 ₇
PT+KPD-A	34.7 ₉	28.9 ₅	25.7 ₄	24.1 ₁₂	19.8 ₂	18.5 ₁₀	26.8 ₄	23.5 ₁₀	19.2 ₁	18.0 ₁₀
FT	4.6 ₇	4.6 ₇	11.5 ₄₁	11.5 ₄₁	9.5 ₁₁	9.5 ₁₁	6.1 ₁₇	6.1 ₁₇	8.5 ₂₄	8.5 ₂₄
PT; FT	20.6 ₁₄	20.6 ₁₄	29.2 ₂₀	29.2 ₂₀	24.9 ₁₁	24.9 ₁₁	23.2 ₁₆	23.2 ₁₆	24.7 ₁₅	24.7 ₁₅
PT+KPD-A; FT	23.3 ₅	23.3 ₅	32.5 ₁₀	32.3 ₉	27.1 ₉	27.1 ₈	26.7 ₁₁	26.7 ₁₂	27.5 ₂	27.4 ₂
PT+KPD-R; FT	26.4 ₁₉	26.4 ₁₉	37.0₁₀	36.6₈	30.6₁₀	30.3₈	28.9₁₄	28.5₁₂	30.9₅	30.6₄

Table 8: Absent and present keyphrase performance using Greedy Search for Transformer One2Set models in a semi-supervised setup. KPD represents KPDrop. PT represents pre-training on the synthetic data (UC). PT+KPD-R or PT+KPD-A represents pre-training on UC with KPD-R or KPD-A respectively. FT represents fine tuning or training on the low resource labelled data (LR). PT (or PT+KPD-A or PT+KPD-R) followed by “; FT” represents that the pre-training was followed by fine-tuning of the pre-trained model on LR. We bold the overall best scores. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1).

Models	Inspec		NUS		Krapivin		SemEval		KP20k	
	F@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5	F1@M	F1@5
Transformer One2Set (Beam Search)										
Absent Keyphrase Performance										
PT	0.2 ₀	0.7 ₂	0.2 ₀	0.6 ₃	0.1 ₀	0.2 ₁	0.2 ₀	0.2 ₂	0.1 ₀	0.3 ₀
PT+KPD-R	0.4₀	1.4 ₂	0.5 ₀	1.4 ₂	0.4 ₀	1.2 ₁₂	0.5 ₁	1.2 ₂	0.3 ₀	1.2 ₁
PT+KPD-A	0.3 ₀	0.8 ₁	0.3 ₀	0.4 ₀	0.2 ₀	0.3 ₁	0.2 ₀	0.3 ₁	0.1 ₀	0.4 ₀
FT	0.2 ₀	0.3 ₂	0.6 ₀	0.8 ₂	0.3 ₀	1.0 ₂	0.4 ₁	0.8 ₁	0.3 ₀	0.9 ₀
PT; FT	0.3 ₀	0.8 ₁	0.5 ₀	1.6 ₃	0.3 ₀	1.4 ₁	0.4 ₀	0.6 ₁	0.3 ₀	1.1 ₁
PT+KPD-A; FT	0.3 ₀	1.0 ₃	0.5 ₀	1.8 ₄	0.4 ₀	1.7 ₃	0.4 ₀	1.5 ₂	0.3 ₀	1.7 ₀
PT+KPD-R; FT	0.4₀	1.7₁	1.0₀	5.4₄	0.7₀	4.5₂	0.9₀	3.1₂	0.5₀	3.6₀
Present Keyphrase Performance										
PT	19.3 ₈	28.6 ₁₅	16.2 ₃	22.7 ₉	10.3 ₃	17.5 ₄	18.0 ₉	23.4 ₁₁	10.5 ₄	17.0 ₅
PT+KPD-R	22.5₄	31.9₁₅	16.9 ₉	23.6 ₇	11.3₅	18.0 ₆	19.6₁₀	22.9 ₁₅	10.8 ₅	17.6 ₇
PT+KPD-A	20.7 ₃	28.8 ₈	17.2₄	23.5 ₉	11.1 ₂	17.9 ₁₀	19.0 ₅	22.9 ₁₅	11.3₂	17.8 ₉
FT	12.6 ₁₀	11.3 ₅	15.0 ₁₅	23.7 ₉	9.3 ₉	14.2 ₇	12.3 ₁₀	15.3 ₁₁	10.9 ₈	16.8 ₁
PT; FT	17.4 ₂	24.3 ₈	14.9 ₃	32.1 ₆	9.5 ₁₂	25.0 ₂	14.8 ₅	27.4 ₁₁	10.5 ₂	25.6 ₃
PT+KPD-A; FT	17.7 ₃	25.9 ₂	15.7 ₁	33.2 ₁₁	10.1 ₂	25.8 ₅	15.8 ₂	26.7 ₁₂	10.9 ₁	26.8 ₂
PT+KPD-R; FT	20.4 ₄	31.9₉	15.4 ₂	37.1₅	10.3 ₂	29.2₇	16.6 ₂	30.4₁₈	10.2 ₂	29.8₁

Table 9: Absent and present keyphrase performance using Beam Search for Transformer One2Set models in a semi-supervised setup. KPD represents KPDrop. PT represents pre-training on the synthetic data (UC). PT+KPD-R or PT+KPD-A represents pre-training on UC with KPD-R or KPD-A respectively. FT represents fine tuning or training on the low resource labelled data (LR). PT (or PT+KPD-A or PT+KPD-R) followed by “; FT” represents that the pre-training was followed by fine-tuning of the pre-trained model on LR. We bold the overall best scores. Subscripts represent standard deviation (e.g., 31.1₁ represents 31.1 ± 0.1).

Models	Inspec		Krapivin		SemEval		KP20k	
	F1@M	F1@5C	F1@M	F1@5C	F1@M	F1@5C	F1@M	F1@5C
ExHiRD Greedy	2.2	1.1	4.3	2.2	2.5	1.7	3.2	1.6
ExHiRD Greedy+KPD-R	3.5	2.0	6.8	3.7	5.1	3.7	5.3	2.7
T5 Greedy	2.5	1.4	5.3	2.8	2.3	1.6	3.6	1.8
T5 Greedy+KPD-R	3.2	1.8	7.0	4.2	3.8	2.9	5.7	3.1

Table 10: Absent keyphrase performance of ExHiRD and T5. We bold the best scores per block.