

Unsupervised Keyphrase Extraction via Interpretable Neural Networks

Anonymous ACL submission

Abstract

Keyphrase extraction aims at automatically extracting a list of “important” phrases representing the key concepts in a document. Prior approaches for unsupervised keyphrase extraction resort to heuristic notions of phrase importance via embedding similarities or graph centrality, requiring extensive domain expertise to develop them. Our work presents an alternative operational definition: phrases that are most useful for predicting the topic of a text are keyphrases. To this end, we propose INSPECT—a self-explaining neural framework for identifying influential keyphrases by measuring the predictive impact of input phrases on the downstream task of topic classification. We show that this novel approach not only alleviates the need for ad-hoc heuristics but also achieves state-of-the-art results in unsupervised keyphrase extraction in 3 out of 4 diverse datasets across two domains: scientific publications and news articles. Ultimately, our study suggests a new usage of interpretable neural networks as an intrinsic component in NLP systems, and not only as a tool for explaining model predictions to humans.

1 Introduction

Keyphrase extraction is crucial for processing and understanding long documents in specialized (e.g., scientific, medical) domains (Mekala and Shang, 2020; Betti et al., 2020; Wang et al., 2019). The task is challenging, as the notion of phrase importance is context- and domain-dependent. For example, scientific terminology has key importance in processing scientific documents (Bekhuis, 2015; Gábor et al., 2016), whereas fine-grained entities and events are important in news summarization (Pighin et al., 2014; Balachandran et al., 2021; Li et al., 2016). Therefore, developing general keyphrase annotation guidelines and curating representative hand-labeled datasets is not feasible. This motivates the need for generalizable unsupervised approaches to keyphrase extraction.

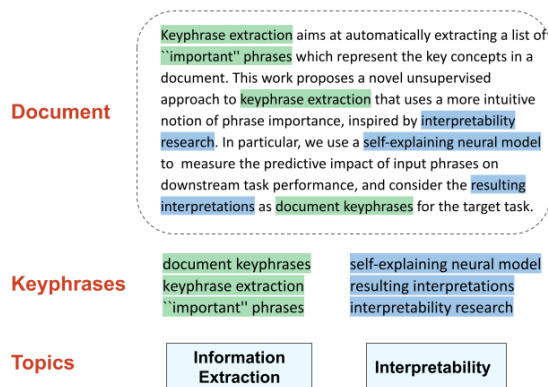


Figure 1: A comprehensive set of keyphrases should highlight important phrases for all major topics in a document. INSPECT presents a method, that leverages interpretable neural models to identify such latent keyphrase useful for predicting topics in a document.

So far, unsupervised approaches to keyphrase extraction have primarily relied on heuristic notions of phrase importance (Mihalcea and Tarau, 2004; Shang et al., 2018; Campos et al., 2018). Popular proxies for phrase importance include phrase clustering based on statistical features like word density (Florescu and Caragea, 2017a; Campos et al., 2018) and structural features like graph centrality (Bougouin et al., 2013). However, such approaches do not yield high-quality keyphrases in new domains as they require domain experts to carefully construct appropriate heuristics (Mani et al., 2020).

In this work, we present an alternative approach. We measure the importance of phrases in a document by their influence on classifying the topics in a text from a set of pre-defined topics. These topics are extracted unsupervisedly in a pre-processing stage. Pre-neural methods have used keyphrases to help identify topics in a document (Wallach, 2006; Wang et al., 2008; Liu and Yang, 2009). We hypothesize that end-to-end neural models also latently use keyphrases to represent documents and perform downstream tasks. Consequently, if we can interpret model decisions via highlighting salient

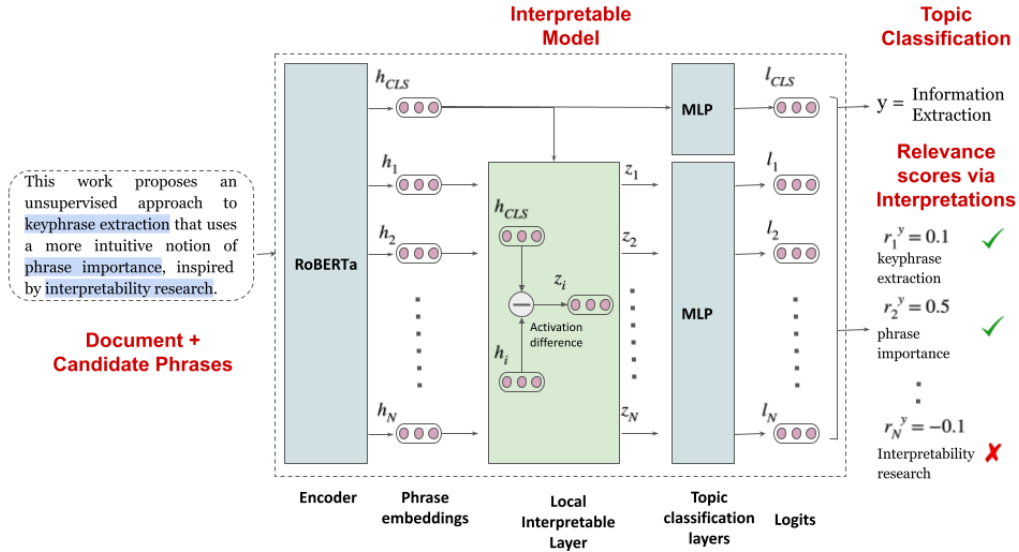


Figure 2: Overview of INSPECT . We first extract candidate phrases and their representations using RoBERTa. The representations of the input without the contribution of each phrase is constructed and provided to the topic classifier. During inference, the model predicts the topics in the document and compares the logits to compute importance scores for each phrase, where a higher score signifies more influence on prediction.

and influential features (phrases) used for prediction, we can identify such keyphrases. Inspired by this intuition, we propose INSPECT—a novel framework that uses interpretable text classifiers to highlight keyphrases important for predicting the topics in a text. Specifically, we adapt an interpretable classifier SelfExplain (Rajagopal et al., 2021) to jointly predict the topic of an input document and to identify the salient phrases influencing the prediction. We consider the resulting interpretations as keyphrases for the input document (§2).

Through extensive experiments, we show that INSPECT is generalizable and can be easily adapted to new domains. We present the applicability of INSPECT in the scientific and news domain (§3). Results on four benchmark datasets show that INSPECT improves keyphrase extraction performance over baselines by up to 11.4% F1 (§4), outperforming the state-of-the-art in unsupervised keyphrase extraction on 3 out of 4 datasets. Importantly, INSPECT alleviates the need for heuristics and expert-labelled annotations, and thus can be applied to a wide range of domains and problems where keyphrase extraction is important. Our results confirm that the latent keyphrases obtained from an interpretable model correlate with human annotated keyphrases, opening new avenues for research on interpretable models for information extraction.¹

¹Code and data will be publicly released.

2 The INSPECT Framework

The goal of the INSPECT framework is to extract important keyphrases in long documents. Following a hypothesis that neural text classifiers latently leverage important keyphrases for predicting topics in text, INSPECT extracts keyphrases through interpreting the classification decisions. It builds upon an interpretable model, SelfExplain (Rajagopal et al., 2021), which learns to attribute text classification decisions to relevant phrases in the input. However, SelfExplain was designed and tested in supervised settings and for single-sentence classification; in this work we explore its extension to unsupervised keyphrase extraction from long documents. In what follows, we describe the base SelfExplain model (§2.1) and the distant supervision setup for *topic classification* (§2.2). We outline the training mechanism to jointly predict topics and highlight salient phrases in the document as model interpretations (§2.3) and finally extract the resulting phrase interpretations as important keyphrases in the document (§2.4). The framework overview is shown in Figure 2.

2.1 Base Interpretable Model

Feature attribution methods for model interpretability include two predominant approaches, (i) post-hoc interpretations of a trained model (Jin et al., 2020; Kennedy et al., 2020; Lundberg and Lee, 2017; Ribeiro et al., 2016), and (ii) intrinsically (by-

design) interpretable models (Alvarez-Melis and Jaakkola, 2018; Rajagopal et al., 2021). We adopt the latter approach, specifically SelfExplain (Rajagopal et al., 2021) as our phrase attribution model, as the model directly produces interpretations, though in principal any phrase based interpretability techniques could be employed.

SelfExplain augments a pre-trained transformer-based model (RoBERTa (Liu et al., 2019) in our case) with a local interpretability layer (LIL) and a global interpretability layer (GIL) which are trained to produce local (relevant features from input sample) and global (relevant samples from training data) interpretations respectively. The model can be trained for any text classification tasks using gold task supervision, and produces local and global interpretations along with model predictions. Since our goal is to identify important phrases from the input sample, we use only the LIL layer. The LIL layer takes as input a sentence x and a set of candidate phrases $CP^x = cp_1, cp_2, \dots, cp_N$ and quantifies the contribution of a particular phrase for prediction through the activation difference (Shrikumar et al., 2017; Montavon et al., 2017) between the phrase and sentence representations.

2.2 Distant Supervision via Topic Prediction

Obtaining annotations for keyphrases in specialized domains is challenging for supervised keyphrase extraction (Mani et al., 2020). Instead, we train the interpretable model in a distant supervision setup for multi-class topic classification and use model interpretations to identify keyphrases, without any keyphrase annotations. Topical information about a document are known to be essential for identifying diverse keyphrases (Bougouin et al., 2013; Sterckx et al., 2015). Further, a comprehensive set of keyphrases should represent the various major topics in the document to be useful for different long document applications (Liu et al., 2010). We hypothesize that by using topic classification as our end-task, our model will learn to highlight—via interpretations it is designed to provide—important and diverse keyphrases in the input document.

While certain domains like news articles have extensive datasets with human annotated topic labels, others like scientific articles or legal documents require significant effort for human annotation. INSPECT can be trained using annotated topic labels when they exist. In other domains where such annotations are scarce, INSPECT can be trained using

labels extracted unsupervisedly using topic models (Gallagher et al., 2017). Experiments in §4 show results using both settings.

2.3 Keyphrase Relevance Model

SelfExplain is designed to process single sentences and uses all the phrases spanning non-terminals in a constituency parser as units (candidate phrases) for interpretation. This is computationally expensive for our use-case. To facilitate long document topic classification, we instead define the set of noun phrases (NPs) as the interpretable units, which aligns with prior work in keyphrase extraction of using noun phrases as initial candidate phrases (Shang et al., 2018; Mihalcea and Tarau, 2004; Bougouin et al., 2013). INSPECT splits a long document into constituent passages, extracts NPs as candidates, and attributes the contribution of each NP for predicting the topics covered in the passage.

For each text block x in the input document, we preprocess and identify a set of candidate phrases $CP^x = cp_1, cp_2, \dots, cp_N$ where N is the number candidate phrases in x . From the base RoBERTa model, we obtain contextual [CLS] representations of the entire text block h_{CLS} and individual tokens. We compute phrase representations $h_1 \dots h_N$ for each candidate in $CP^x = \{cp_1, cp_2, \dots, cp_N\}$ by taking the sum of the RoBERTa representations of each token in the phrase cp_i .

To compute the relevance of each phrase, we construct a representation of the input without the contribution of the phrase, z_i , using the activation differences between the two representations. We then pass it to a classifier layer in the local interpretability module to obtain the label distribution for prediction.

$$z_i = g(h_i) - g(h_{CLS}); \quad l_i = f(W^T z_i + b) \quad (1)$$

where g is the ReLU activation function and W and b are the weights and bias of the classifier. Here l_i denotes the label distribution obtained on passing the phrase-level representations z_i through a classification layer f which is either the sigmoid or the softmax function depending on the prediction task (multi-label versus multi-class). We denote the label distribution from the base RoBERTa model for predicting the output using the whole input block as l_{CLS} . We train the model using the cross entropy loss \mathcal{L}_y with respect to the multi-label gold topics y_t and an explanation specific loss \mathcal{L}_e using the mean of all phrase-level label distributions such

that $l_e = \text{mean}(l_i)$.

$$\mathcal{L}_y = - \sum_{t=1}^T y_t \log(l_{CLS}), \quad \mathcal{L}_e = - \sum_{t=1}^T y_t \log(l_e) \quad (2)$$

The classifier is regularized jointly with α parameter using explanation and classification loss:

$$\mathcal{L} = (1 - \alpha)\mathcal{L}_y + \alpha\mathcal{L}_e,$$

2.4 Inference

During inference, for each predicted label $y \in Y$, where Y denotes set of all predicted labels for input text x , INSPECT calculates an importance score r_i^y with respect to the predicted label y using the difference between the label distribution l_i^y for a candidate phrase c_i and the one obtained using the entire input l_{CLS}^y as $r_i^y = l_{CLS}^y - l_i^y$.

This score denotes the influence of a candidate keyphrase on the predicted topic. A higher score is caused by a high shift in label distribution when using the representation of the input without the contribution of the phrase, indicating that the phrase is highly relevant for predicting the topic. Since the relevance scores are computed with respect to a particular predicted topic and its label distribution, the scores for the same input are not comparable across different predicted topics in multi-label classification (since label distributions can vary in magnitude). To aggregate important keyphrases across all predicted topics, we pick the ones that positively impact prediction for each topic (having a positive influence score) as a set of keyphrases.

$$KP(x) = [CP_i \mid \forall r_i^y > 0; y \in Y; i \in \{1 : N\}]$$

3 Experimental Setup

3.1 Evaluation Datasets

We evaluate INSPECT in two domains using four popular keyphrase extraction datasets—scientific publications (SemEval-2017 (Augenstein et al., 2017a), SciERC (Luan et al., 2018), SciREX (Jain et al., 2020)) and news articles (500N-KPCrowd (Marujo et al., 2013)). Dataset details and statistics are listed in Table 5 in the Appendix §A.1.

3.2 Topic Labels

We create distant supervision for INSPECT by labeling the above datasets using document topics as labels. We leverage existing topic annotations when such annotations exist. In the 500N-KPCrowd news based dataset, we use existing topic labels

(tags or categories such as Sports, Politics, Entertainment) in a one-class classification setting. For the scientific publications domain, we use topic models (Gallagher et al., 2017) to extract $T = 75$ topics where each document can be labeled with multiple topics. The scientific domain datasets are trained in a multi-label classification setup.

3.3 Training Data and Settings

We evaluate the generalizability of INSPECT in two experimental settings:

- INSPECT:** In this setting, for each of our datasets we train the model for topic prediction using only the documents and topic labels from the training set of the dataset. We then evaluate using the held-out test data from the dataset. The training data in this setting, is most closely aligned to the test data, as the documents are of the same topic distribution.
- INSPECT-ZeroShot:** Here, the model is trained using a large external dataset of documents and topic labels from a similar domain. The model is then evaluated on the test data of each dataset. The training data here is of a similar domain (e.g. ICLR papers for scientific domain), but not necessarily of similar topic distribution as the test data (e.g. SemEval-2017 has Physics papers which have different topics than ICLR papers). In this setting, we use data from ICLR OpenReview² papers for scientific domain and BBC News articles for news domain. We collect over 8,317 full papers from ICLR and obtained 75 topic labels using topic modeling³. We removed 22 topic labels that were uninformative (list in Appendix Table 6) and used the rest to train our model in a multi-label classification setup. The BBC News corpus (Greene and Cunningham, 2006) consists of 2,225 news article documents, each annotated with one of five topics (business, entertainment, politics, sport, or tech).

We pre-process each document by splitting it into text blocks of size 512 tokens, where consecutive blocks overlap with a stride size of 128. Following Shang et al. (2018), for each block we consider all Noun Phrases (NPs) as candidate phrases

²<https://openreview.net/group?id=ICLR>.
cc

³https://github.com/gregversteeg/corex_topic

Dataset	Method	F1 Score		
		Micro	Macro	Weighted
SciERC	RoBERTa	0.842	0.651	0.767
	INSPECT	0.836	0.658	0.771
SciREX	RoBERTa	0.609	0.404	0.641
	INSPECT	0.628	0.442	0.697
SemEval17	RoBERTa	0.819	0.613	0.731
	INSPECT	0.822	0.611	0.744
500N-KPCrowd	RoBERTa	0.916	0.880	0.910
	INSPECT	0.938	0.904	0.939
ICLR	RoBERTa	0.729	0.456	0.699
	INSPECT	0.743	0.492	0.733
BBC News	RoBERTa	0.880	0.851	0.876
	INSPECT	0.902	0.886	0.894

Table 1: Proxy Task (Topic prediction) performance. Our INSPECT method outperforms a strong RoBERTa baseline on Micro, Macro and Weighted F1 scores.

and extract them using a Noun Phrase extractor from the Berkeley Neural Parser⁴. All hyperparameters were chosen based on development set performance on SciERC. Our final models were trained with a batch size of 8 a learning rate of 2e-5 for 10 epochs. The classification layer dimension was 64 and α was 0.5.⁵

3.4 Baselines

We compare our method against seven unsupervised keyphrase extraction techniques — Yake (Campos et al., 2018), TF-IDF (Florescu and Caragea, 2017a), TopicRank (Bougouin et al., 2013) AutoPhrase (Shang et al., 2018; Liu et al., 2015), SifRank (Sun et al., 2020), AttentionRank (Ding and Luo, 2021) and UAE-CCRank (Liang et al., 2021). Out of the chosen baselines, Yake, TF-IDF and AutoPhrase are statistical, TopicRank is graph-based and SifRank, AttentionRank and UAE-CCRank are neural embedding based methods. Following prior work and task guidelines (Augenstein et al., 2017a; Jain et al., 2020), INSPECT produces **span level** keyphrases and distinguishes each occurrence of a keyphrase. In contrast, methods like SifRank, AttentionRank, and UAE-CCRank are phrase level keyphrase extractors which don’t provide span level outputs. To maintain common evaluation, we adapt these methods to span level keyphrase extraction by matching each output keyphrase to all occurrences of the phrase in the document. As our method applies a cutoff on relevance scores and picks any phrase with a positive relevance score as a keyphrase, we

⁴<https://pypi.org/project/benepar/>

⁵Details on our hyperparameter search is shared in the appendix §A.2

cannot be directly compared with baselines which rank candidate phrases and pick top-K phrases as important. To establish a fair setting for evaluation, we use the average of the number of keyphrase predictions from our model as the ‘K’ across all baselines.

3.5 Evaluation Metrics

Topic Prediction Evaluation: To ensure high-quality interpretations from our model, it is imperative that it performs well on topic prediction. We first evaluate INSPECT’s performance on topic prediction using average F1 scores across all labels.

Keyphrase Extraction Evaluation: For our primary evaluation of keyphrase extraction, we evaluate using span match of our predictions and the true labels (human annotated keyphrases). Prior works (Shang et al., 2018; El-Beltagy and Rafea, 2009; Bougouin et al., 2013) have mainly focused on *exact match* performance. However, a recent survey highlights that the measure is highly restrictive (Papagiannopoulou and Tsoumakas, 2019) as simple variations in preprocessing can misalign phrases giving an inaccurate representation of the model’s capabilities (Boudin et al., 2016).

Alternatively, *partial span match* using the word level overlap between the predicted and gold span ranges, has also been explored (Rousseau and Vazirgiannis, 2015). But, it is sometimes lenient in scoring. Papagiannopoulou and Tsoumakas (2019) suggest *average of the exact and partial matching* as an appropriate metric based on empirical studies. Therefore, we evaluate performance using the average of the exact and partial match F1 scores between predicted and true phrases keyphrases.

4 Results

4.1 Topic Prediction with INSPECT

First, we compare INSPECT’s effectiveness in classifying the topics with the corresponding non-interpretable encoder baseline, using micro, macro, and weighted F1 score of the classifier’s predictions compared to gold standard annotations. The results in Table 1 show that our approach outperforms a strong RoBERTa (Liu et al., 2019) baseline for topic prediction across all of our evaluation datasets. The difference is more pronounced in larger datasets (SciREX, ICLR, and BBC News), and strong performance on the topic classification task provides confidence that highlighted interpretations are for relevant and major topics in the text.

Dataset	Method	Exact Match F1	Partial Match F1	Avg Exact Partial F1
SciERC	TF-IDF	0.0627	0.2860	0.1743
	TopicRank	0.2533	0.5680	0.4110
	Yake	0.2230	0.5125	0.3678
	AutoPhrase	0.0961	0.3145	0.2053
	AttentionRank	0.3461	0.4690	0.4075
	UKE CCRank	0.3584	0.4804	0.4194
	INSPECT	0.3108	0.5524	0.4316
SciREX	TF-IDF	0.1521	0.3690	0.2605
	TopicRank	0.2298	0.4122	0.3210
	Yake	0.1840	0.3734	0.2787
	AutoPhrase	0.1814	0.4236	0.3025
	AttentionRank	0.2554	0.2198	0.2376
	UKE CCRank	0.0419	0.0759	0.0589
	INSPECT	0.2397	0.4127	0.3262
SemEval17	TF-IDF	0.0610	0.2698	0.1654
	TopicRank	0.2240	0.4312	0.3276
	Yake	0.1687	0.3644	0.2665
	AutoPhrase	0.0790	0.3404	0.2097
	AttentionRank	0.2408	0.3442	0.2925
	UKE CCRank	0.2427	0.345	0.2938
	INSPECT	0.2594	0.5185	0.3889
500N-KPCrowd	TF-IDF	0.1034	0.3520	0.2277
	TopicRank	0.1060	0.2346	0.1703
	Yake	0.1380	0.3551	0.2465
	AutoPhrase	0.1590	0.3608	0.2599
	AttentionRank	0.3032	0.3442	0.3237
	UKE CCRank	0.1729	0.2873	0.2303
	INSPECT	0.1608	0.3920	0.2764

Table 2: Span-match results for unsupervised keyphrase extraction across datasets in the INSPECT setting. Best performance is indicated in Bold. **Our model outperforms baselines on average of exact and partial F1 scores.**

4.2 Keyphrase Span Match Performance

Next, we study the utility of INSPECT in highlighting keyphrases via model interpretations. The results for INSPECT are detailed in Table 2 and, for INSPECT-ZeroShot in Table 3.⁶

Results in Table 2 show that even with access to only training set of documents from each dataset, on 3 out of 4 datasets INSPECT outperforms all baselines with ~ 3.24 average F1 improvements. In the news domain (500-KPCrowd dataset) INSPECT has low exact match scores but higher partial match scores indicating misalignments between predicted and gold spans. Additionally, 500N-KPCrowd annotates all instances of a keyphrase as a reference span which favours phrase level methods like AttentionRank in the current evaluation setup. In SciREX, we observe very poor performance of UKE CCRank as it ranks common phrases like “image”, “label”, “method”, etc, very high.

In the INSPECT-ZeroShot setting, with access to a larger dataset of external documents, our model outperforms prior methods in 3 out of 4 datasets

⁶As SifRank uses external data to augment the model, we compare it with INSPECT-ZeroShot for a fair comparison. Baselines that don’t make use of any external corpus are only included in the INSPECT setting evaluation.

with ~ 3.2 points average F1 improvements. In the 500N-KPCrowd dataset, INSPECT performs comparably to SifRank with improved Partial Match F1. As Table 3 illustrates, we notice that the model consistently performs better in the INSPECT-ZeroShot setting when compared with the INSPECT setting, showing that the method benefits from more training data. Our results further show that variations in topic distribution between training and test data don’t significantly impact results. INSPECT can thus benefit from large unlabeled documents from similar domains to improve results.

INSPECT improves performance in settings with human annotated topics (news) as well as when topics are extracted using unsupervised topic modeling (scientific). Additionally, most baselines rely on carefully constructed pre- and post-processing to eliminate common phrases and produce high-quality candidates (Liang et al., 2021; Ding and Luo, 2021; Sun et al., 2020). In contrast, INSPECT achieves competitive results without domain expertise and processing for extracting quality keyphrases. Therefore, INSPECT can be easily adapted to new domains without human annotations for topics and with minimal domain knowledge.

Dataset	Method	Exact Match F1	Partial Match F1	Avg Exact Partial F1
SciERC	TF-IDF	0.2162	0.4434	0.3298
	AutoPhrase	0.2416	0.6130	0.4273
	SifRank	0.2248	0.7357	0.4803
	Our	0.4371	0.7114	0.5743
SciREX	TF-IDF	0.1780	0.4008	0.2894
	AutoPhrase	0.2583	0.4993	0.3788
	SifRank	0.1234	0.3957	0.2595
	Our	0.2601	0.4893	0.3747
SemEval17	TF-IDF	0.1810	0.3398	0.2604
	AutoPhrase	0.1104	0.4874	0.2989
	SifRank	0.2804	0.6336	0.457
	Our	0.3246	0.6218	0.4732
500N-KPCrowd	TF-IDF	0.1398	0.3578	0.2488
	AutoPhrase	0.1701	0.3918	0.2805
	SifRank	0.1847	0.4125	0.2986
	Our	0.1776	0.4194	0.2985

Table 3: Span-match results for unsupervised keyphrase extraction in INSPECT-ZeroShot (trained on ICLR and BBC News corpus). Best performance is indicated in Bold. **INSPECT outperforms most baselines.**

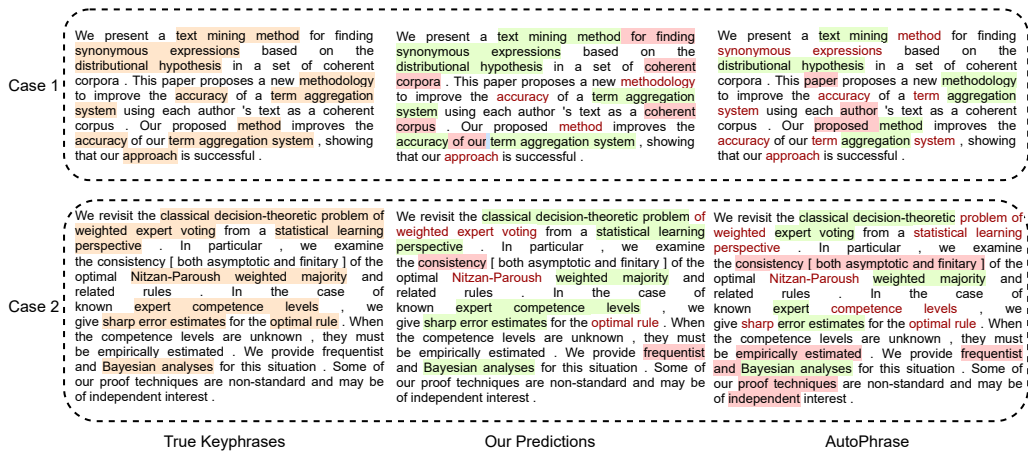


Figure 3: Two data points randomly chosen from the SciERC dataset. Orange spans represent gold standard annotations. Green spans in the predictions represent correctly predicted spans, whereas red spans are spans wrongly predicted as being keyphrases and red text are keyphrases that the model did not identify.

Our results demonstrate that phrase attribution techniques from interpretability literature can be leveraged to identify high-quality document keyphrases by measuring predictive impact of input phrases on topic prediction. Crucially, as these keyphrases correlate with human annotated keyphrases, our results validate our initial hypothesis that neural models latently use document keyphrases for tasks like topic classification.

5 Discussion

Here, we present an analysis on the common error types in INSPECT and discuss the strengths and weaknesses of INSPECT using qualitative examples.

Entity Type Analysis: We leverage the entity type information in SciERC to observe the performance of INSPECT on specific types of keyphrases. From Table 4, we see that INSPECT performs best

on keyphrases labelled as *Scientific Terms* and *Materials*. *Generic* phrases and *Metrics* are usually not representative of topical content, and thus, our method performs poorly on them. On manual analysis, we noticed that many phrases marked as *Task* are very unique and infrequent, making them harder to identify. A high partial match recall but a low exact match recall for *Method* type suggest that many predicted keyphrases are misaligned with the gold labels. We believe that alternative downstream tasks can be explored in future to help tailor our approach to capture specific types of entities, based on application requirements.

Qualitative Analysis In Figure 3 we show two randomly selected abstracts from the SciERC dataset. We see that INSPECT tends to extract longer phrases compared to AutoPhrase, which tends to extract mostly unigrams or bigrams. Since

Type	Recall	
	Exact	Partial
Metric	60.65	78.34
Task	58.27	90.45
Material	72.17	86.69
Scientific Term	78.87	95.13
Method	65.31	95.41
Generic	63.16	86.06

Table 4: Exact and partial span match recall scores for different types of keyphrases on the SciERC dataset.

noun phrases can overlap, we observe that our model sometimes predicts overlapping phrases. Overall, our approach is able to extract more relevant phrases than the baseline. Both INSPECT and AutoPhrase tend to miss generic phrases like ‘approach’ (e.g., as seen in case 1). We hypothesize that since the downstream task in INSPECT is to identify topics, it would lead the model to focus on phrases more relevant for detecting the topic of the document. Also, this could make INSPECT miss highly targeted phrases (which usually consist of proper nouns) like *Nitzan-Paroush* in case 2.

INSPECT tends to extract longer compound phrases connected by functional words. Potentially, post-processing to remove overlapping and compound phrases might lead to even higher performance on datasets with smaller keyphrases. Case 2 in Figure 3 also demonstrates the ability of predicting complete phrases, like ‘classical decision-theoretic problem’, instead of AutoPhrase’s prediction – ‘classical decision-theoretic’ which is incomplete.

6 Related Work

Unsupervised keyphrase extraction is typically treated as a ranking problem, given a set of candidate phrases (Shang et al., 2018; Campos et al., 2018; Florescu and Caragea, 2017a). A standard pipeline (1) extracts candidate phrases; (2) scores phrase relevance; and (3) ranks the phrases based on their scores. Broadly, prior approaches can be categorized as statistical, graph-based, embedding-based, or language model based methods; Papiagiannopoulou and Tsoumakas (2019) provide a detailed survey.

Statistical methods exploit notions of information theory directly. Common approaches include TF-IDF based scoring (Florescu and Caragea, 2017a) of phrases with other co-occurrence statistics to enhance performance (Liu et al., 2009; El-Beltagy and Rafea, 2009). Campos et al. (2018)

shows the importance of incorporating statistical information of the context of each phrase to improve performance. Statistical approaches typically treat different instances of a phrase equally, which is a limitation.

Graph-based techniques, on the other hand, broadly aim to form a graph of candidate phrases connected based on similarity to each other. Then core components of the graph are chosen as key phrases. Amongst these, PageRank (Brin and Page, 1998) and TextRank (Mihalcea and Tarau, 2004) assign scores to nodes based on their influence. A common extension is to use weights on the edges denoting the strength of connection (Wan and Xiao, 2008; Rose et al., 2010; Bougouin et al., 2013). Position Rank (Florescu and Caragea, 2017b) and SGRank (Danesh et al., 2015) combine the ideas from statistical, word co-occurrence and positional information. Some approaches, especially applied in the scientific document setting, make use of citation graphs (Gollapalli and Caragea, 2014; Wan and Xiao, 2008), and external knowledge bases (Yu and Ng, 2018) to improve keyphrase extraction. In this work, we focus our approach on a general unsupervised keyphrase extraction setting applicable to any domain where such external resources may not be present.

Finally, embedding based techniques (Bennani-Smires et al., 2018; Papiagiannopoulou and Tsoumakas, 2018) make use of word-document similarity using word embeddings (Sun et al., 2020; Liang et al., 2021), while language-model based techniques use the uncertainty when predicting words to decide informativeness (Tomokiyo and Hurst, 2003). Ding and Luo (2021) uses attention scores to calculate phrase importance with the document in an unsupervised manner.

7 Conclusion and Future Work

In this work, we introduced INSPECT, a novel approach to unsupervised keyphrase extraction. Our framework uses a neural model that explains text classification decisions to extract keyphrases via phrase-level feature attribution. Using four standard datasets in two domains, we show that INSPECT outperforms prior methods and establishes state-of-art results in 3 out of 4 datasets.. Through qualitative and quantitative analysis, we show that INSPECT can produce high-quality and relevant keyphrases. INSPECT presents applications of interpretable models beyond explanations for humans.

References

David Alvarez-Melis and Tommi S Jaakkola. 2018. Towards robust interpretability with self-explaining neural networks. *Neurips*.

Isabelle Augenstein, Mrinal Das, Sebastian Riedel, Lakshmi Vikraman, and Andrew McCallum. 2017a. **SemEval 2017 task 10: ScienceIE - extracting keyphrases and relations from scientific publications**. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 546–555, Vancouver, Canada. Association for Computational Linguistics.

Isabelle Augenstein, Mrinal Das, Sebastian Riedel, Lakshmi Vikraman, and Andrew McCallum. 2017b. **SemEval 2017 task 10: ScienceIE - extracting keyphrases and relations from scientific publications**. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 546–555, Vancouver, Canada. Association for Computational Linguistics.

Vidhisha Balachandran, Artidoro Pagnoni, Jay Yoon Lee, Dheeraj Rajagopal, J. Carbonell, and Yulia Tsvetkov. 2021. Structsum: Summarization via structured representations. In *EACL*.

Tanja Bekhuis. 2015. Keywords, discoverability, and impact. *Journal of the Medical Library Association : JMLA*, 103 3:119–20.

Kamil Bennani-Smires, Claudiu Musat, Andreea Hossmann, Michael Baeriswyl, and Martin Jaggi. 2018. Simple unsupervised keyphrase extraction using sentence embeddings. *arXiv preprint arXiv:1801.04470*.

Arianna Betti, Martin Reynaert, Thijs Ossenkople, Yvette Oortwijn, Andrew Salway, and Jelke Bloem. 2020. **Expert concept-modeling ground truth construction for word embeddings evaluation in concept-focused domains**. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6690–6702, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Florian Boudin, Hugo Mougard, and Damien Cram. 2016. How document pre-processing affects keyphrase extraction performance. In *NUT@COLING*.

Adrien Bougouin, Florian Boudin, and Béatrice Daille. 2013. **TopicRank: Graph-based topic ranking for keyphrase extraction**. In *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, pages 543–551, Nagoya, Japan. Asian Federation of Natural Language Processing.

Sergey Brin and Lawrence Page. 1998. **The anatomy of a large-scale hypertextual web search engine**. *Computer Networks*, 30:107–117.

Ricardo Campos, Vítor Mangaravite, Arian Pasquali, Alípio Mário Jorge, Célia Nunes, and Adam Jatowt. 2018. A text feature based automatic keyword extraction method for single documents. In *European conference on information retrieval*, pages 684–691. Springer. 621
622
623
624
625
626

Soheil Danesh, Tamara Sumner, and James H. Martin. 2015. **SGRank: Combining statistical and graphical methods to improve the state of the art in unsupervised keyphrase extraction**. In *Proceedings of the Fourth Joint Conference on Lexical and Computational Semantics*, pages 117–126, Denver, Colorado. Association for Computational Linguistics. 627
628
629
630
631
632
633

Haoran Ding and Xiao Luo. 2021. **AttentionRank: Unsupervised keyphrase extraction using self and cross attentions**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1919–1928, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 634
635
636
637
638
639
640

Samhaa R. El-Beltagy and Ahmed Rafea. 2009. **Kp-miner: A keyphrase extraction system for english and arabic documents**. *Information Systems*, 34(1):132–144. 641
642
643
644

Corina Florescu and Cornelia Caragea. 2017a. A new scheme for scoring phrases in unsupervised keyphrase extraction. In *European Conference on Information Retrieval*, pages 477–483. Springer. 645
646
647
648

Corina Florescu and Cornelia Caragea. 2017b. **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. 649
650
651
652
653
654

K. Gábor, Haïfa Zargayouna, D. Buscaldi, I. Tellier, and Thierry Charnois. 2016. Semantic annotation of the acl anthology corpus for the automatic analysis of scientific literature. In *LREC*. 655
656
657
658

Ryan J Gallagher, Kyle Reing, David Kale, and Greg Ver Steeg. 2017. Anchored correlation explanation: Topic modeling with minimal domain knowledge. *Transactions of the Association for Computational Linguistics*, 5:529–542. 659
660
661
662
663

Sujatha Das Gollapalli and Cornelia Caragea. 2014. Extracting keyphrases from research papers using citation networks. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, AAAI’14*, page 1629–1635. AAAI Press. 664
665
666
667
668

Derek Greene and Pádraig Cunningham. 2006. Practical solutions to the problem of diagonal dominance in kernel document clustering. In *Proc. 23rd International Conference on Machine Learning (ICML’06)*, pages 377–384. ACM Press. 669
670
671
672
673

Sarthak Jain, Madeleine van Zuylen, Hannaneh Hajishirzi, and Iz Beltagy. 2020. **Scirex: A challenge dataset for document-level information extraction**. 674
675
676

677					
678					
679	Xisen Jin, Zhongyu Wei, Junyi Du, Xiangyang Xue,				
680	and Xiang Ren. 2020. Towards hierarchical impor-				
681	tance attribution: Explaining compositional seman-				
682	tics for neural sequence models. In <i>International</i>				
683	<i>Conference on Learning Representations.</i>				
684	Brendan Kennedy, Xisen Jin, Aida Mostafazadeh Da-				
685	vani, Morteza Dehghani, and Xiang Ren. 2020. Con-				
686	textualizing hate speech classifiers with post-hoc ex-				
687	planation. In <i>Proceedings of the 58th Annual Meet-</i>				
688	<i>ing of the Association for Computational Linguistics,</i>				
689	pages 5435–5442, Online. Association for Computa-				
690	tional Linguistics.				
691	Wei Li, Lei He, and Hai Zhuge. 2016. Abstrac-				
692	tive news summarization based on event semantic				
693	link network. In <i>Proceedings of COLING 2016,</i>				
694	<i>the 26th International Conference on Computational</i>				
695	<i>Linguistics: Technical Papers,</i> pages 236–246, Os-				
696	aka, Japan. The COLING 2016 Organizing Commit-				
697	tee.				
698	Xinnian Liang, Shuangzhi Wu, Mu Li, and Zhoujun Li.				
699	2021. Unsupervised keyphrase extraction by jointly				
700	modeling local and global context. In <i>Proceedings</i>				
701	<i>of the 2021 Conference on Empirical Methods in</i>				
702	<i>Natural Language Processing,</i> pages 155–164, On-				
703	line and Punta Cana, Dominican Republic. Associa-				
704	tion for Computational Linguistics.				
705	Jialu Liu, Jingbo Shang, Chi Wang, Xiang Ren, and Ji-				
706	awei Han. 2015. Mining quality phrases from mas-				
707	sive text corpora. In <i>Proceedings of the 2015 ACM</i>				
708	<i>SIGMOD International Conference on Management</i>				
709	<i>of Data,</i> pages 1729–1744.				
710	Nan Liu and Christopher C. Yang. 2009. Keyphrase				
711	extraction for labeling a website topic hierarchy. In				
712	<i>ICEC.</i>				
713	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-				
714	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,				
715	Luke Zettlemoyer, and Veselin Stoyanov. 2019.				
716	Roberta: A robustly optimized bert pretraining ap-				
717	proach. <i>arXiv preprint arXiv:1907.11692.</i>				
718	Zhiyuan Liu, Wenyi Huang, Yabin Zheng, and				
719	Maosong Sun. 2010. Automatic keyphrase extrac-				
720	tion via topic decomposition. In <i>EMNLP.</i>				
721	Zhiyuan Liu, Peng Li, Yabin Zheng, and Maosong				
722	Sun. 2009. Clustering to find exemplar terms for				
723	keyphrase extraction. In <i>Proceedings of the 2009</i>				
724	<i>conference on empirical methods in natural lan-</i>				
725	<i>guage processing,</i> pages 257–266.				
726	Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh				
727	Hajishirzi. 2018. Multi-task identification of enti-				
728	ties, relations, and coreference for scientific knowl-				
729	edge graph construction. In <i>Proceedings of the 2018</i>				
730	<i>Conference on Empirical Methods in Natural Lan-</i>				
731	<i>guage Processing,</i> pages 3219–3232, Brussels, Bel-				
732	gium. Association for Computational Linguistics.				
	Scott M. Lundberg and Su-In Lee. 2017. A unified				733
	approach to interpreting model predictions. In <i>Pro-</i>				734
	<i>ceedings of the 31st International Conference on</i>				735
	<i>Neural Information Processing Systems, NIPS’17,</i>				736
	page 4768–4777, Red Hook, NY, USA. Curran As-				737
	sociates Inc.				738
	Kaushik Mani, Xiang Yue, Bernal Jimenez Gutierrez,				739
	Yungui Huang, Simon Lin, and Huan Sun. 2020.				740
	Clinical phrase mining with language models. In				741
	<i>2020 IEEE International Conference on Bioinfor-</i>				742
	<i>matics and Biomedicine (BIBM),</i> pages 1087–1090.				743
	IEEE.				744
	Luis Marujo, Márcio Viveiros, and João Paulo da Silva				745
	Neto. 2013. Keyphrase cloud generation of broad-				746
	cast news. In <i>Proceeding of Interspeech 2011: 12th</i>				747
	<i>Annual Conference of the International Speech Com-</i>				748
	<i>munication Association.</i>				749
	Dheeraj Mekala and Jingbo Shang. 2020. Contextual-				750
	ized weak supervision for text classification. In <i>Pro-</i>				751
	<i>ceedings of the 58th Annual Meeting of the Associa-</i>				752
	<i>tion for Computational Linguistics,</i> pages 323–333,				753
	Online. Association for Computational Linguistics.				754
	Rada Mihalcea and Paul Tarau. 2004. Textrank: Bring-				755
	ing order into text. In <i>Proceedings of the 2004 con-</i>				756
	<i>ference on empirical methods in natural language</i>				757
	<i>processing,</i> pages 404–411.				758
	Grégoire Montavon, Sebastian Lapuschkin, Alexander				759
	Binder, Wojciech Samek, and Klaus-Robert Müller.				760
	2017. Explaining nonlinear classification decisions				761
	with deep taylor decomposition. <i>Pattern Recogni-</i>				762
	<i>tion,</i> 65:211–222.				763
	Eirini Papagiannopoulou and Grigorios Tsoumakas.				764
	2018. Local word vectors guiding keyphrase ex-				765
	traction. <i>Information Processing & Management,</i>				766
	54(6):888–902.				767
	Eirini Papagiannopoulou and Grigorios Tsoumakas.				768
	2019. A review of keyphrase extraction. <i>CoRR,</i>				769
	abs/1905.05044.				770
	Daniele Pighin, M. Cornolti, Enrique Alfonseca, and				771
	Katja Filippova. 2014. Modelling events through				772
	memory-based, open-ie patterns for abstractive sum-				773
	marization. In <i>ACL.</i>				774
	Dheeraj Rajagopal, Vidhisha Balachandran, E. Hovy,				775
	and Yulia Tsvetkov. 2021. Selfexplain: A self-				776
	explaining architecture for neural text classifiers.				777
	<i>ArXiv,</i> abs/2103.12279.				778
	Marco Tulio Ribeiro, Sameer Singh, and Carlos				779
	Guestrin. 2016. "why should i trust you?": Explain-				780
	ing the predictions of any classifier. <i>Proceedings of</i>				781
	<i>the 22nd ACM SIGKDD International Conference</i>				782
	<i>on Knowledge Discovery and Data Mining.</i>				783
	Stuart Rose, Dave Engel, Nick Cramer, and Wendy				784
	Cowley. 2010. Automatic keyword extraction from				785
	individual documents. <i>Text mining: applications</i>				786
	<i>and theory,</i> 1:1–20.				787

788	François Rousseau and Michalis Vazirgiannis. 2015.	Yang Yu and Vincent Ng. 2018. Wikirank: Improving	843
789	Main core retention on graph-of-words for single-	keyphrase extraction based on background knowl-	844
790	document keyword extraction. In <i>European Con-</i>	edge. <i>arXiv preprint arXiv:1803.09000</i> .	845
791	ference on Information Retrieval, pages 382–393.		
792	Springer.		
793	Jingbo Shang, Jialu Liu, Meng Jiang, Xiang Ren,	A Appendix	846
794	Clare R Voss, and Jiawei Han. 2018. Automated	A.1 Evaluation Datasets	847
795	phrase mining from massive text corpora. <i>IEEE</i>	SemEval-2017 (Augenstein et al., 2017a) consists	848
796	<i>Transactions on Knowledge and Data Engineering</i> ,	of 500 abstracts taken from 12 AI conferences cov-	849
797	30(10):1825–1837.	ering Computer Science, Material Science, and	850
798	Avanti Shrikumar, Peyton Greenside, and Anshul Kun-	Physics. The entities are annotated with Process,	851
799	daje. 2017. Learning important features through	Task, and Material labels, which form the funda-	852
800	propagating activation differences. In <i>International</i>	mental concepts in scientific literature. Identifica-	853
801	<i>Conference on Machine Learning</i> , pages 3145–3153.	tion of the keyphrases was subtask A of the Scien-	854
802	PMLR.	ceIE SemEval task (Augenstein et al., 2017b).	855
803	Lucas Sterckx, Thomas Demeester, Johannes Deleu,	SciERC (Luan et al., 2018) extends SemEval-	856
804	and Chris Develder. 2015. When topic models	2017 by annotating more entity types, relations,	857
805	disagree: Keyphrase extraction with multiple topic	and co-reference clusters to include broader cover-	858
806	models. <i>Proceedings of the 24th International Con-</i>	age of general AI. The dataset was annotated by a	859
807	ference on World Wide Web.	single domain expert who had high (76.9%) agree-	860
808	Yi Sun, Hangping Qiu, Yu Zheng, Zhongwei Wang,	ment with three other expert annotators on 12%	861
809	and Chaoran Zhang. 2020. Sifrank: A new base-	subset of the dataset.	862
810	line for unsupervised keyphrase extraction based on	SciREX (Jain et al., 2020) is a document-level	863
811	pre-trained language model. <i>IEEE Access</i> , 8:10896–	information extraction dataset, covering entity iden-	864
812	10906.	tification and n-ary relation formation using salient	865
813	Takashi Tomokiyo and Matthew Hurst. 2003. A lan-	entities. Human and automatic annotations were	866
814	guage model approach to keyphrase extraction. In	used to annotate 438 full papers with salient enti-	867
815	<i>Proceedings of the ACL 2003 Workshop on Multi-</i>	tities, with a distant supervision from the Papers	868
816	<i>word Expressions: Analysis, Acquisition and Treat-</i>	With Code ⁷ corpus. This dataset can help verify	869
817	<i>ment - Volume 18, MWE '03</i> , page 33–40, USA. As-	the performance of models on full papers.	870
818	sociation for Computational Linguistics.	500N-KPCrowd (Marujo et al., 2013) is a	871
819	Hanna M. Wallach. 2006. Topic modeling: beyond	keyphrase extraction dataset in the news domain.	872
820	bag-of-words. <i>Proceedings of the 23rd international</i>	This data consists of 500 articles from 10 topics	873
821	<i>conference on Machine learning</i> .	annotated by multiple Amazon Mechanical Turk	874
822	Xiaojun Wan and Jianguo Xiao. 2008. Single doc-	workers for important keywords. Following the	875
823	ument keyphrase extraction using neighborhood	baselines on this datasets, we pick keywords that	876
824	knowledge. In <i>AAAI</i> , volume 8, pages 855–860.	were among the top two most frequently chosen by	877
825	Canhui Wang, Min Zhang, Liyun Ru, and Shaop-	the human annotators. Since no span-level infor-	878
826	ing Ma. 2008. An automatic online news topic	mation for these keywords is given, we annotate all	879
827	keyphrase extraction system. <i>2008 IEEE/WIC/ACM</i>	occurrences of the chosen keywords in the docu-	880
828	<i>International Conference on Web Intelligence and</i>	ment to obtain a list of span labels, which we use	881
829	<i>Intelligent Agent Technology</i> , 1:214–219.	to evaluate all the models.	882
830	Wenbo Wang, Yang Gao, He-Yan Huang, and Yuxiang	A.2 Implementation Details	883
831	Zhou. 2019. Concept pointer network for abstrac-	Here, we present the hyper-parameters for all exper-	884
832	tive summarization. In <i>Proceedings of the 2019 Con-</i>	iments along with their corresponding search space.	885
833	ference on Empirical Methods in Natural Language	We chose all hyperparameters based on the devel-	886
834	Processing and the 9th International Joint Confer-	opment set performance on the SciERC dataset.	887
835	ence on Natural Language Processing (EMNLP-	We considered RoBERTa (Liu et al., 2019) and	888
836	IJCNLP), pages 3076–3085.	XL-NET (Yang et al., 2019) based encoders and	889
837	Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-	finally chose RoBERTa for faster compute times.	890
838	bonell, Russ R Salakhutdinov, and Quoc V Le. 2019.		
839	Xlnet: Generalized autoregressive pretraining for		
840	language understanding. In <i>Advances in Neural In-</i>		
841	formation Processing Systems, volume 32. Curran		
842	Associates, Inc.		

⁷<https://paperswithcode.com/>

Dataset	Type	Split	Total docs	Avg words per doc	Avg keyphrases per doc
SciERC	Scientific	Train	350	130	16
		Dev	50	130	16
		Test	100	134	17
SciREX	Scientific	Train	306	5601	353
		Dev	66	5484	354
		Test	66	6231	387
SemEval17	Scientific	Train	350	160	21
		Dev	50	193	27
		Test	100	186	23
500N-KPCrowd	News	Train	400	430	193
		Dev	50	465	86
		Test	50	420	116
BBC News	News	All	2225	385	-
ICLR	Scientific	All	8317	6505	-

Table 5: Description about the datasets. Average words and keyphrases per document are rounded to the nearest whole number. ICLR and BBC News are used in INSPECT-ZeroShot setting for training and don't have any labelled keyphrase data.

S.No.	Top words from removed topic
1	proposed;propose novel;propose;proposed method;method
2	generalization;study;analysis;suggest;provide
3	outperforms;existing;existing methods;outperforms stateofheart;methods
4	state;art;state art;shortterm;current state
5	effectiveness;demonstrate effectiveness;source;effectiveness proposed;student
6	training;training data;training set;training process;model training
7	experimental;experimental results;results;results demonstrate;experimental results demonstrate
8	experiments;extensive;extensive experiments;experiments demonstrate;conduct
9	performance;improves;significantly;improve;improved
10	recent;shown;recent work;recent advances;success
11	achieves;introduce;competitive;achieves stateofheart;introduce new
12	trained;model trained;models trained;networks trained;trained using
13	present;paper present;present novel;work present;monte
14	widely;parameters;widely used;proposes;paper proposes
15	simple;benchmark datasets;benchmark;propose simple;simple effective
16	prior;approach;sampling;continuous;prior work
17	program;introduces;programs;future;paper introduces
18	solve;challenging;able;complex;challenging problem
19	challenge;current;challenges;open;current stateofheart
20	rate;good;good performance;l;regime
21	works;previous works;existing works;focus;scenarios
22	evaluate;evaluation;tackle;tackle problem;evaluate method

Table 6: 22 Generic topics removed from the 75 topic labels learned using topic modeling on ICLR data.

891 We experimented with learning-rates from the set
892 of $1e-5, 2e-5, 5e-5, 1e-4$ and $2e-4$. We chose $2e-5$
893 as the final learning rate. Our batch size of 8 was
894 chosen after experimenting with 4, 8, 12 and 16.
895 The size of the weights matrix in the classification
896 layer was chosen to be 64 from a set of 16, 32, 64
897 and 128. The α parameter used for regularization
898 was fixed at 0.5. We tried values between 0.1 and
899 0.9 and did not find significant difference. We saved
900 the model based on best weighted F1 on the topic
901 prediction task. All training runs took less than
902 3 hours on 2 Nvidia 2080Ti GPUs, except on the
903 ICLR dataset, which took 8 hours. All results are
904 from a single run.