

# Domain Representative Keywords Selection: A Probabilistic Approach

Pritom Saha Akash<sup>1</sup> Jie Huang<sup>1</sup> Kevin Chen-Chuan Chang<sup>1</sup>  
Yunyaoli Li<sup>2\*</sup> Lucian Popa<sup>3</sup> ChengXiang Zhai<sup>1</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign, USA

<sup>2</sup>Apple, USA

<sup>3</sup>IBM Research, USA

{pakash2, jeffhj, kcchang, czhai}@illinois.edu  
yunyaoli@apple.com, lpopa@us.ibm.com

## Abstract

We propose a probabilistic approach to select a subset of a *target domain representative keywords* from a candidate set, contrasting with a context domain. Such a task is crucial for many downstream tasks in natural language processing. To contrast the target domain and the context domain, we adapt the *two-component mixture model* concept to generate a distribution of candidate keywords. It provides more importance to the *distinctive* keywords of the target domain than common keywords contrasting with the context domain. To support the *representativeness* of the selected keywords towards the target domain, we introduce an *optimization algorithm* for selecting the subset from the generated candidate distribution. We have shown that the optimization algorithm can be efficiently implemented with a near-optimal approximation guarantee. Finally, extensive experiments on multiple domains demonstrate the superiority of our approach over other baselines for the tasks of keyword summary generation and trending keywords selection.<sup>1</sup>

## 1 Introduction

*Domain representative keywords* are the core knowledge of a *target domain* of interest. A target domain can be a broad area of science like *computer science (CS)* or its sub-field *artificial intelligence (AI)*. Acquiring domain representative keywords benefits various natural language processing (NLP) tasks such as information summarization, organization, and extraction. For instance, acquiring a set of domain representative keywords is an important first step in organizing domain knowledge with a taxonomy of keywords (Zhang et al., 2018). Moreover, tagging documents (Chen et al., 2017) with domain representative keywords helps

\* This work was done while the author was at IBM Research, USA.

<sup>1</sup>Code and data are available at <https://github.com/pritomsaha/keyword-selection>

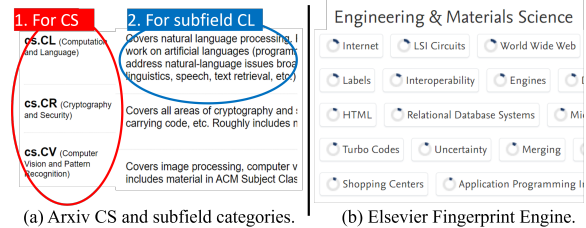


Figure 1: Screenshots of example applications.

to facilitate search or recommendation in a domain. For another example, summarizing a domain using its trending keywords for a specific time frame helps researchers get a snapshot of research trends or emerging areas of interest, e.g., new emerging security vulnerabilities.

In reality, while representing a domain, the desired keywords often depend on a given *context domain*. E.g., if we are interested in representing the CS domain with the context of general knowledge (all areas of knowledge), keywords like *model*, *data*, *information* make sense in distinguishing CS from general knowledge. However, if the context is general science, those keywords are not *distinctive* enough to distinguish CS from other areas like mathematics or physics. Instead, the keywords *machine learning*, *data mining*, *deep learning* make more sense in this case. Therefore, it is important to contrast with a known context domain while representing a target domain, but most of the existing work ignored this. An application of this is shown in Fig. 1 (a) from the *arxiv category taxonomy*<sup>2</sup>. We can see a shift of categories from general to more specific subcategories of research areas: CS → CL. Knowing CS categories (partially shown in (1)) as the context, the keywords specified in (2) (i.e., *speech*, *text retrieval*) are more appropriate to represent CL than keywords overlapped with other CS categories.

Moreover, the number of keywords that need to

<sup>2</sup>[https://arxiv.org/category\\_taxonomy](https://arxiv.org/category_taxonomy)

be selected depends on the nature of the applications. E.g., users are interested in a quick, high-level overview with fewer keywords while summarizing a particular domain. On the other hand, while building a controlled vocabulary representative of a certain domain, the number of keywords is naturally large. An application of the controlled vocabulary is illustrated in Figure 1 (b). It shows the fingerprint visualization of a CS researcher generated by Elsevier Fingerprint Engine<sup>3</sup>, a system for research profiling. The researcher's profile is summarized using keywords in *Engineering & Material Science* domain. However, we can see that some non-representative keywords like *labels* and *merging* are used in the summarization. Therefore, having a representative controlled vocabulary for each domain will facilitate this application for expressively representing a researcher profile.

We thus propose the **problem of domain representative keywords selection**. As input, we are given a set of candidate keywords, a target and a context domain represented by their corresponding corpora, and a size  $k$ . As output, we aim to select a subset consisting of  $k$  keywords from the given candidate set such that the subset best represents the target domain contrasting with the context domain. Here, we assume that the candidate keywords are from the target domain and can be implicitly extracted from the given target domain corpus or externally given keywords for that domain.

From the above problem and discussion, we have identified that the solution for the problem needs to meet the following two **requirements**: (1) the selected keywords should be distinctive to the target domain contrastive with a context domain; (2) the selected keywords should represent the target domain as a whole within the specific size constraint. None of the existing work satisfies all of them. Previously, research has been conducted on automatic keyword extraction (Hätty et al., 2017; Meng et al., 2017; Alzaidy et al., 2019; Wang et al., 2020) and phrase mining (Liu et al., 2015; Shang et al., 2018). However, their main focus is to extract terms from single/multiple documents without considering whether the extracted terms are distinctive to a target domain contrastive with a context. There is also some previous research (Liu et al., 2015; Shang et al., 2018; Lu et al., 2019; Huang et al., 2021) that tries to find fine-grained domain-

specific keywords from the text. However, these approaches mostly rank keywords based on their specificity to a corpus (or domain) rather than selecting a predefined number of keywords with a global objective of representing the target domain. Therefore, in this work, we propose a solution to satisfy all the specified requirements.

The first **challenge** on fulfilling the requirements is *contrasting the target and context domains*. Among candidate keywords, the distinctive keywords may have similar corpus statistics (i.e., frequency from target domain corpus) with many non-distinctive popular keywords. Therefore, simply filtering out highly frequent keywords may lose many distinctive keywords for a target domain. Instead, it is more intuitive to say that the keywords that frequently appear in both target and context corpora are often not distinctive keywords for the target domain. It inspires us to leverage the two-component mixture model (MM) (Zhai and Lafferty, 2001) concept to generate the candidate keywords distribution contrasting with the context domain. As far as we know, this is the first work to utilize a mixture model mechanism for keywords selection.

The second **challenge** is the *representation under a size constraint*. If we simply select the top distinctive keywords based on the MM-generated distribution, we may end up with redundant keywords that may fall short in representing the target domain as a whole. Hence, it is more intuitive to consider selecting keywords with a domain representation objective. Therefore, we cast this as an optimization problem of selecting  $k$  keywords that *coarsen* the candidate distribution adapting the concept of *statistical machine translation* (Brown et al., 1993) with the objective of minimizing the divergence between the initial and coarsened distributions of candidate keywords.

In summary, as our **contributions** in this paper, **firstly**, we propose a new problem formulation named *domain representative keywords selection*. **Secondly**, we propose a framework for solving the problem consisting of two steps: (1) generating candidate keywords distribution using a two-component mixture model mechanism and (2) selecting a subset of keywords utilizing the generated distribution with an introduced optimization algorithm. **Thirdly**, we prove that our proposed optimization problem can be efficiently solved with a near-optimal approximation ratio. **Finally**, to validate the effectiveness of our approaches, we con-

<sup>3</sup><https://www.elsevier.com/solutions/elsevier-fingerprint-engine>

duct extensive experiments on multiple domains for different tasks demonstrating the superiority of our framework against strongly designed baselines.

## 2 Related Work

The problem of *domain representative keywords selection* is related to the automatic keyword extraction (AKE) problem. AKE focuses on extracting or generating the most prominent keywords from single/multiple documents. Existing methods for AKE can be classified into two categories: supervised and unsupervised keyword extraction. Early supervised methods consider AKE as a binary classification problem (Witten et al., 1999; Turney, 2000) by learning a classifier from annotated documents to predict whether a candidate phrase is a keyword or not. Recently, *deep learning* has been used for the supervised AKE. E.g., (Meng et al., 2017) uses an encoder-decoder-based framework to generate keywords where (Alzaidy et al., 2019) addresses AKE as a sequence labeling problem. Unsupervised AKE methods mostly apply graph-based ranking mechanisms utilizing semantic relatedness measure between keywords (Mihalcea and Tarau, 2004). Besides, linguistic (Handler et al., 2016) and semantic (Bennani-Smires et al., 2018) approaches have also been used for unsupervised AKE.

However, the main focus of the above studies is to describe single/multiple documents rather than domain-specific keywords extraction. To solve this problem, several researches (Liu et al., 2015; Shang et al., 2018; Lu et al., 2019; Wang et al., 2020; Huang et al., 2021) have been conducted on domain-specific fine-grained keyword extraction. E.g., (Huang et al., 2021) propose an algorithm for measuring the relevance of a keyword in a particular domain. However, this approach requires a user to provide some seed domain-relevant terms for supervising the algorithm. Moreover, the above approaches only consider ranking keywords based on their domain specificity (or relevance). None of them deals with the problem of domain representative keyword selection with a specific size constraint.

The mixture model used for generating keywords distribution in our approach is related to the research on probabilistic topic models (Hofmann, 2001; Blei et al., 2003) and comparative text mining (Sarawagi et al., 2003; Zhai et al., 2004). However, the difference between our approach and

these studies is that rather than finding multiple latent topics or themes from a collection or multiple collections of documents, we model a target domain corpus as a distribution of unigram language model contrastive with a context model.

## 3 Proposed Methodology

Our proposed framework consists of two steps: (1) generating distribution for the candidate keywords and (2) selecting a subset that best represents the target domain utilizing the generated distribution.

### 3.1 Keywords Distribution Generation

To select keywords, how do we represent a target domain in contrast with a context one? One naive solution can be the frequency distribution of keywords in the target corpus. However, this distribution is biased towards common but possibly non-distinctive keywords (e.g., *data*, *method* and *model* in CS), which may not differentiate the target (e.g., CS) from the context (e.g., Physics) domain. On the other hand, among candidate keywords, the distinctive keywords may have similar corpus statistics (i.e., frequency from target domain corpus) with many non-distinctive common keywords. Therefore, it is not easy to separate those desired target domain keywords from non-distinctive common keywords using simple statistics calculated from the target domain corpus. E.g., keyword *algorithm* is more distinctive than *method*, but both are popular keywords in CS domain. Therefore, simply filtering out highly frequent keywords may lose many distinctive keywords for a target domain.

To handle the above problem, we regard the target corpus as a mixture of two unigram language models. Specifically, the corpus is assumed to be generated from a mixture of two multinomial component models. One model is the known background model  $\theta_B$  (computed from the context corpus), which models the non-distinctive common keywords in the target and context corpora. The other one is the target domain model ( $\theta_D$ ) that needs to be estimated and concerned for prioritizing distinctive keywords in that domain.

Formally, let  $\mathcal{C}$  be the target domain corpus from which we are interested to find the keyword distribution, then the log-likelihood value (LLV) of generating  $\mathcal{C}$  from this mixture model is

$$\begin{aligned} \log p(\mathcal{C}|\theta_D) = \\ \sum_{t_i \in V} c(t_i, \mathcal{C}) \log[(1 - \lambda)p(t_i|\theta_D) + \lambda p(t_i|\theta_B)], \end{aligned} \quad (1)$$

where  $V$  is the candidate keywords set and  $c(t_i, \mathcal{C})$  is the frequency of keyword  $t_i$  in  $\mathcal{C}$ .  $\lambda$  refers to the mixing weight of the  $\theta_B$ . In other words,  $\lambda$  controls the amount of “background noise” in the corpus we want to be modeled by  $\theta_B$ . We assume  $\theta_B$  and  $\lambda$  to be known, and  $\theta_D$  be estimated. Specifically,  $\theta_B$  is the probability distribution calculated from the context domain corpus.

In principle, we can estimate  $\theta_D$  using any optimization methods. E.g., the Expectation-Maximization (EM) algorithm (Dempster et al., 1977) is one of them and can be used to compute a maximum likelihood estimate with the following updating formulas:

$$\begin{aligned} p^{(n)}(z = 0|t_i) &= \frac{(1 - \lambda)p^{(n)}(t_i|\theta_D)}{(1 - \lambda)p^{(n)}(t_i|\theta_D) + \lambda p^{(n)}(t_i|\theta_B)}, \\ p^{(n+1)}(t_i|\theta_D) &= \frac{c(t_i, \mathcal{C})p^{(n)}(z = 0|t_i)}{\sum_{t_j \in V} c(t_j, \mathcal{C})p^{(n)}(z = 0|t_j)}, \end{aligned}$$

where  $p(z = 0|t_i)$  refers how likely  $t_i$  is from  $\theta_D$ . The estimated  $\{p(t_i|\theta_D) \cdots p(t_N|\theta_D)\}$  is used as candidate keywords distribution.

### 3.2 Keyword Subset Selection

After acquiring a distribution of candidate keywords, we find a subset with a size  $k$  to represent the target domain. One possible solution is to select top  $k$  keywords based on the candidate distribution ( $\theta_D$ ) generated by the *mixture model* (MM). Hence, the keywords with high distinctiveness to the target domain contrasting with the context domain will be selected. However, one problem with this approach is that the selected keywords may fall short in representing the target domain by only selecting some redundant distinctive keywords.

To solve the above problem, we view the subset selection as a *distribution coarsening* problem. Specifically, we want to use a subset to estimate the candidate distribution (i.e., coarsened distribution). As defined in the previous section, a domain is a distribution of keywords (i.e., candidate distribution). Therefore, for a subset of keywords to represent the domain, the coarsened distribution by the subset should closely approximate the candidate distribution of that domain.

Formally, let  $\mathcal{P} = \{p(t_i) \cdots p(t_N)\}$  be the candidate distribution, we compute a coarsened distribution  $\tilde{\mathcal{P}} = \{\tilde{p}(t_i) \cdots \tilde{p}(t_N)\}$  by subset  $S$  and  $\tilde{p}(t_i)$  for each  $t_i \in V$  is calculated as:

$$\tilde{p}(t_i) = \sum_{t_j \in S} p(t_i|t_j)p(t_j), \quad (2)$$

where  $p(t_i|t_j)$  refers to the probability of *semantically translating*  $t_j$  into  $t_i$ . This idea of estimating the probability of each keyword from candidates by a subset is adapted from the statistical machine translation from the same language used in information retrieval (Berger and Lafferty, 1999).

Now to find the subset, we introduce an optimization problem with objective of selecting a subset ( $S$ ) with size  $k$  from candidates ( $V$ ) that minimizes the difference between the LLV of generating  $\mathcal{C}$  by  $\mathcal{P}$  and  $\tilde{\mathcal{P}}$ , respectively. We know that the LLV of generating  $\mathcal{C}$  by  $\mathcal{P}$  is  $\log p(\mathcal{C}) = \sum_{t_i \in V} c(t_i, \mathcal{C}) \log p(t_i)$  where  $c(t_i, \mathcal{C})$  is the frequency of  $t_i$  in  $\mathcal{C}$ . Similarly, the LLV of generating  $\mathcal{C}$  from  $\tilde{\mathcal{P}}$  is  $\log \tilde{p}(\mathcal{C})$ . Hence, given  $|S| = k$ , our optimization objective is:

$$\begin{aligned} S &= \arg \min_{S \subseteq V} \|\log p(\mathcal{C}) - \log \tilde{p}(\mathcal{C})\| \\ &= \arg \min_{S \subseteq V} \left\| \sum_{t_i \in V} c(t_i, \mathcal{C}) \log \frac{p(t_i)}{\tilde{p}(t_i)} \right\| \\ &= \arg \min_{S \subseteq V} \left\| \sum_{t_i \in V} p(t_i) \log \frac{p(t_i)}{\tilde{p}(t_i)} \right\| \\ &= \arg \min_{S \subseteq V} D_{KL}(\mathcal{P} \|\tilde{\mathcal{P}}) = \arg \min_{S \subseteq V} \phi(S), \end{aligned} \quad (3)$$

where  $\phi(S)$  is our objective function, and  $D_{KL}(\mathcal{P} \|\tilde{\mathcal{P}})$  is Kullback–Leibler (KL) divergence (Kullback and Leibler, 1951) between  $\mathcal{P}$  and  $\tilde{\mathcal{P}}$ .

From (2), one obvious question is how to calculate  $p(t_i|t_j)$ . For this, we use *mutual information* (MI) to estimate  $p(t_i|t_j)$  inspired from (Karimzadehgan and Zhai, 2010) where MI is used to estimate a similar model for information retrieval. MI is a good measure to judge relatedness between two terms. In our model, for any two terms  $t_i$  and  $t_j$ , we first compute MI ( $I(t_i; t_j)$ ) between them and normalize it into a probability as below:

$$I(t_i; t_j) = \sum_{b_i, b_j} p(b_i, b_j) \log \frac{p(b_i, b_j)}{p(b_i)p(b_j)}, \quad (4)$$

$$p(t_i|t_j) \approx p_{MI}(t_i|t_j) = \frac{I(t_i; t_j)}{\sum_{t'_j \in V} I(t_i; t'_j)}, \quad (5)$$

where  $b_i$  is a binary variable indicating the presence/absence of  $t_i$ . E.g.,  $p(b_i = 1)$  indicates the ratio of documents containing  $t_i$  and  $p(b_i = 1, b_j = 1)$  indicates the ratio of documents where both  $t_i$  and  $t_j$  co-occur. Here,  $p_{MI}(t_i|t_j)$  gives us the probability of how  $t_j$  relates to  $t_i$ ; intuitively, this probability would be higher when these two terms frequently co-occur in the same document in the target corpus.

**Optimization.** We are interested in finding a subset  $S$  with size  $k$  from  $V$  such that  $\phi(S)$  is minimized, i.e.,  $\arg \min_{S \in V} \phi(S)$  s.t.  $|S| = k$ . This is referred as the *cardinality-constrained optimization* and proven to be NP-hard (Feige, 1998). However, if the objective function  $\phi(S)$  is monotone and submodular, a simple greedy algorithm is guaranteed to obtain an approximation of  $1 - \frac{1}{e}$ . We call a non-negative real valued function  $F$  (to be maximized) *submodular* if it has the property of *diminishing returns* that is  $F(X \cup \{v\}) - F(\{v\}) \geq F(Y \cup \{v\}) - F(\{v\})$  for all  $v \in V$  and  $X \subseteq Y \subseteq V$ . Moreover,  $F$  is said to be *monotone* if  $F(X) \leq F(Y)$  for all  $X \subseteq Y$ .

**Theorem 1.** *For minimizing the objective function  $\phi(\cdot)$ , a simple greedy algorithm obtains an approximation guarantee of  $1 - \frac{1}{e}$ .*

*Proof.* The proof can be found in Appendix A.  $\square$

So, as per Theorem 1, we can obtain a near optimal solution using a simple greedy algorithm. Initially, we have  $S = \emptyset$ , then iteratively update  $S = S \cup \arg \max_{t \in V \setminus S} \mathcal{G}(t|S)$  until  $|S| = k$  where  $\mathcal{G}(t|S) = \phi(S) - \phi(S \cup t)$  is the gain of adding a new term  $t$  to  $S$ . Thanks to the submodularity property of  $\phi(\cdot)$ , this simple greedy algorithm can further be accelerated by lazy greedy algorithm (Minoux, 1978). More specifically, instead of re-computing  $\mathcal{G}(t_i|S)$ ,  $\forall t_i \in V$  in every step, we use a priority queue of sorted gains  $g(t_i)$ ,  $\forall t_i \in V$ . Starting with  $g(t_i) = -\phi(\{t_i\})$ ,  $\forall t_i \in V$ , the algorithm adds a term  $t_i$  to  $S$  if  $g(t_i) \geq \mathcal{G}(t_i|S)$ , otherwise we update  $g(t_i)$  to  $\mathcal{G}(t_i|S)$  and resort the priority queue. This largely improves the efficiency of the algorithm.

## 4 Experiments

This section evaluates our models from different perspectives: (1) the ability to select representative summary keywords for a target domain; (2) the performance for trending keywords selection task in a domain for different time frames.

### 4.1 Experiment Setup

**Datasets.** In our experiments, to test the generality of the proposed approaches, we use two document collections from two domains for constructing target and context corpora for each of the domains. One is abstracts collections from the *arxiv* repository (version 47)<sup>4</sup>, and the other is a collection of newsgroup documents<sup>5</sup>.

**Candidate Keywords.** In our experiments, we use different sets of candidate keywords. For the CS domain, we collected keywords from two external sources named Springer and Aminer (Tang et al., 2008). The Springer CS keyword list is collected through web scraping from Springer<sup>6</sup> and trimmed to 83K based on frequency  $\geq 5$ . The Aminer keyword list is the collection of keywords assigned by authors in CS research papers, and there are approximately 50K keywords in this list. Alongside keywords from external sources, we also created candidate sets extracted from concerning corpus using AutoPhrase (Shang et al., 2018) tool. All the candidate keywords are lemmatized, and several filtering rules are used. For instance, keywords containing only letters, numbers, hyphens are used; stop and single-letter words are removed.

**Baselines.** We compare our models with the following four baseline keyword selection algorithms.

- **Relative Frequency (RF):** Since a keyword is likely to be domain representative when it frequently appears in a domain corpus, we consider a simple approach that selects the top  $k$  frequent keywords based on the relative frequency calculated from the target corpus.
- **Log-odds (LO):** We adapted a method (Monroe et al., 2008) for keyword selection which was introduced to compare words used by two political parties. Recently, (Hughes et al., 2020) used this method for detecting trending terms in *Cybersecurity* forum discussion. In this baseline, we adapt this method to model keywords as a function of a particular domain or time to compute the likelihood of keywords in that domain or time as log-likelihoods (“log-odds”).
- **Page Rank (PR):** This baseline is a graph-based keyword selection method using PageRank (Mihalcea and Tarau, 2004). We build the graph of candidate keywords where each edge weight denotes how closely two keywords are related.

<sup>4</sup><https://www.kaggle.com/Cornell-University/arxiv>

<sup>5</sup><http://qwone.com/jason/20Newsgroups>

<sup>6</sup><https://www.springer.com/gp>

- **Facility Location (FL) Function:** Facility location function is a representation based subset selection measure (Mirchandani and Francis, 1990) used for finding a representative subset of items. Recently, this measure is used for training-data subset selection (Kaushal et al., 2019). In this paper, we adapt this measure as a baseline for selecting subset from candidate keywords set. Specifically, denoting  $rel(t_i, t_j)$  as the relatedness of two keywords  $t_i$  and  $t_j$ , the objective is to select a subset  $S \in V$  that maximizes FL function  $f(S) = \sum_{t_i \in V} \max_{t_j \in S} rel(t_i, t_j)$ .

**Proposed Models.** We have the following three variants of our proposed framework.

- **KL divergence + RF (KL<sub>rf</sub>):** This model is a simple version of our proposed objective function  $D_{KL}(\mathcal{P} \parallel \tilde{\mathcal{P}})$  defined in (3). In this model,  $\mathcal{P}$  is the relative frequency distribution calculated from the target corpus and  $\tilde{\mathcal{P}}$  is coarsened distribution defined in (2).
- **Mixture Model (MM):** In this proposed model, keywords are ranked based on the estimated distribution for the target domain contrasting with a context domain using the mixture model defined in Section 3.1. Based on the distribution, the top  $k$  keywords are selected.
- **KL Divergence + MM (KL<sub>mm</sub>):** This proposed model is similar to KL<sub>rf</sub>. In KL<sub>mm</sub>, instead of using relative frequency, the mixture model estimated keyword distribution is used as  $\mathcal{P}$  in  $D_{KL}(\mathcal{P} \parallel \tilde{\mathcal{P}})$ .

**Implementation Details.** There are some parameters both in baselines and the proposed models we have to set. E.g., the mixing weight  $\lambda$  for the background model in the mixture model is set to two different values based on the specificity of the target domains. Particularly, when we set  $\lambda$  to a small value, the model favors frequent non-informative terms (i.e., domain-specific stop words). Therefore, the larger values are set for  $\lambda$ . In our experiments, for a broad domain like CS, we set  $\lambda$  to 0.9, and for more specific domains (i.e., AI and subtopics in newsgroup), we set  $\lambda$  to 0.99. The reason for these two different values of  $\lambda$  is that more specific domains demand larger  $\lambda$  for selecting distinctive keywords. For optimizing MM, we use Expectation-Maximization (EM) algorithm (Dempster et al., 1977). Since EM does not guarantee the global maxima, in our experiment, we run the algorithm multiple times with random initialization, and the one with the best MLE is chosen to

$k$	RF	LO	PR	FL	KL <sub>rf</sub>	MM	KL <sub>mm</sub>
10	1.0651	1.1001	1.1035	1.0722	1.0651	2.0981	<b>2.1212</b>
20	1.1440	3.2476	2.2682	1.1345	1.1451	4.2626	<b>4.2972</b>
30	2.2134	4.3965	3.3607	3.2682	3.2875	<b>4.4321</b>	4.4169
40	3.3273	4.4929	4.4902	3.3515	3.3660	4.6000	<b>4.6018</b>
50	3.4530	4.6896	4.5734	3.4505	3.4496	<b>5.6826</b>	<b>5.6826</b>
100	4.7812	8.2382	7.0399	4.7626	4.8761	8.2708	<b>8.2824</b>
200	9.7166	11.1047	8.9403	9.5908	8.7045	11.1082	<b>12.0233</b>
500	18.902	19.1719	18.0464	16.6171	17.7441	<b>19.3221</b>	19.2353

Table 1: Category correspondence results

reduce the chance of getting local maxima. As we use mutual information (MI) based on document co-occurrence statistics in our model (defined in (5)), for the fair comparison, in the baseline FL, we also use MI between two keywords  $t_i$  and  $t_j$  to encode the relatedness between them (i.e.,  $rel(t_i, t_j)$ ). Similarly, MI is used for computing edge weight in the PR method.

## 4.2 Experiment Results

### 4.2.1 Summary Keywords Selection

We conduct both quantitative and qualitative studies to evaluate the ability of proposed models to select domain representative summary keywords. For this purpose, we use the abstracts from the arxiv under CS categories as the target corpus. The context corpus is composed of all abstracts in the arxiv repository.

**Quantitative Evaluation.** We create keyword summaries for the CS domain with varying sizes ( $k$ ) for quantitative evaluation. We collected 52 known category keywords from arxiv categories as CS representative ground keywords to evaluate the ability of selected  $k$  summary keywords to represent the target domain when  $k$  varies. The correspondence between  $k$  selected keywords  $S = \{t_1 \cdots t_k\}$  and  $m$  category keywords  $C = \{c_1 \cdots c_m\}$ ,  $CC(S, C)$  is calculated as the summation of the pairwise normalized mutual information (NMI) (Bouma, 2009) between  $S$  and  $C$  i.e.,  $CC(S, C) = \sum_{i,j} \frac{I(t_i; c_j)}{H(t_i; c_j)}$  where  $I(t_i; c_j)$  is calculated following formula from (4) and  $H(t_i; c_j) = -\sum_{b_i, b_j} p(b_i, b_j) \log p(b_i, b_j)$  is the joint entropy of  $t_i$  and  $c_j$ .

From the results on Table 1, using AutoPhrase extracted candidate keywords, we can see that even though no supervision is used, our methods KL<sub>mm</sub> and MM select keywords that best correspond with the known categories outperforming all the baselines (similar results from two more candidate sets are shown in Appendix B). We ob-

Models	Selected 20 keywords in CS
RF	paper, model, <b>algorithm</b> , datum, result, information, <b>graph</b> , state, high, art, single, order, human, research, general, design, <b>deep learning</b> , <b>semantic</b> , knowledge, <b>neural network</b>
LO	<b>algorithm</b> , art, information, <b>semantic</b> , <b>graph</b> , human, <b>deep learning</b> , paper, datum, <b>neural network</b> , <b>machine learning</b> , real world, research, video, <b>robot</b> , communication, language, <b>security</b> , <b>architecture</b> , knowledge
PR	polynomial, research, channel, paper, energy, <b>graph</b> , datum, model, experimental, information, <b>machine learning</b> , <b>software</b> , binary, english, propose method, function, acoustic, upper, solution, algebraic
FL	art, paper, datum, <b>algorithm</b> , model, result, high, information, <b>graph</b> , channel, research, order, single, human, general, <b>deep learning</b> , design, experimental, solution, knowledge
$KL_{rf}$	model, <b>algorithm</b> , paper, datum, state, <b>graph</b> , result, information, high, art, human, research, design, single, <b>semantic</b> , order, <b>deep learning</b> , energy, general, <b>neural network</b>
MM	<b>algorithm</b> , art, <b>semantic</b> , <b>deep learning</b> , human, <b>neural network</b> , <b>convolutional neural network</b> , <b>machine learning</b> , real world, video, information, <b>robot</b> , research, language, communication, <b>security</b> , <b>architecture</b> , <b>privacy</b> , <b>deep neural network</b> , <b>label</b>
$KL_{mm}$	<b>algorithm</b> , art, <b>semantic</b> , <b>deep learning</b> , human, <b>security</b> , <b>neural network</b> , real world, <b>convolutional neural network</b> , communication, <b>machine learning</b> , <b>robot</b> , language, video, research, <b>privacy</b> , <b>label</b> , information, <b>software</b> , <b>architecture</b>

Keywords distinctive to the CS domain are **highlighted** (annotated by authors).

Table 2: Summary keywords in CS Domain

serve that there is a good improvement of result from MM to  $KL_{mm}$ . However, this is not true for  $KL_{rf}$  and the RF baseline. The reason is that the relative frequencies from the target corpus favor the non-distinctive common keywords (e.g., *model* and *method*). As described in Section 3.2,  $KL_{rf}$  tries to select the subset of keywords that best estimate the original candidate distribution. Hence, it also favors those common keywords to attain the nearest estimation of the original distribution.

On the other hand, the MM-generated distribution assigns larger probabilities to distinctive keywords of the target domain, contrasting with the context domain. Therefore, selecting a keyword subset by  $KL_{mm}$  with close estimation of the MM generated distribution also favors distinctive keywords with the domain representative objective. Furthermore, one interesting observation is that when  $k$  is smaller, the selected keywords by  $KL_{mm}$  tend to summarize the domain better than that of MM. The primary reason for this is that  $KL_{mm}$  prefers to select more non-redundant keywords than MM while  $k$  is smaller, which we later discuss from Table 2.

Models	2000-2009	2010-2019	2020-2021
RF	0.6289	0.6640	0.6493
LO	0.6813	0.7199	0.7238
PR	0.6626	0.6970	0.6826
FL	0.6172	0.6848	0.6528
$KL_{rf}$	0.6282	0.6792	0.6516
MM	<b>0.6908</b>	0.7331	0.7898
$KL_{mm}$	0.6898	<b>0.7763</b>	<b>0.7944</b>

Table 3: Results using trending ground truth keywords

Models	2000-2009	2010-2019	2020-2021
RF	0.3593	0.3195	0.323
LO	0.3956	0.3326	0.4043
PR	0.3705	0.3211	0.3641
FL	0.3591	0.3217	0.3336
$KL_{rf}$	0.3583	0.3189	0.3239
MM	0.4104	0.3468	0.5145
$KL_{mm}$	<b>0.4165</b>	<b>0.3523</b>	<b>0.5215</b>

Table 4: Results generated using Google Trends

**Qualitative Evaluation.** For the qualitative evaluation, we show the summary keywords selected by different algorithms in the CS domain from AutoPhrase extracted candidate keywords in Table 2 (similar additional results are shown in Appendix C). This study aims to observe the difference between the proposed models and baselines in selecting summary keywords. We can see that our models (MM and  $KL_{mm}$ ) outperform all the baselines by selecting the most number of CS representative keywords. We also observe that the LO baseline method also selects a comparable amount of distinctive keywords. The reason is its use of a contrastive method like MM for selecting keywords for a particular corpus compared to a context corpus.

However, our models MM and  $KL_{mm}$  tend to select more representative keywords than the LO method. E.g., we can see that our methods select keywords like *privacy*, *software* and *convolutional neural network* instead of keywords that LO selects like *graph* and *paper*, *data*. Another observation is that the keywords selected by PR are mostly those keywords (i.e., *experimental*, *data* and *function*) that have a broad association with other words. However, these keywords as an unit do not convey much information about the domain.

Now to see the difference between our models MM and  $KL_{mm}$ , we see the difference between their selected keywords. As stated before, we can

	2000-2009	2010-2019	2020-2021
RF	paper, problem, <b>algorithm</b> , model, method, approach, <b>system</b> , information, result, datum, set, application, number, user, word, performance, language, order, time, case	model, method, paper, approach, image, problem, datum, task, algorithm, dataset, performance, network, result, feature, system, training, application, work, number, object	model, method, task, datum, approach, dataset, image, paper, performance, problem, training, algorithm, network, feature, result, system, work, <b>application</b> , <b>deep learning</b> , experiment
LO	<b>logic program</b> , rule, manipulator, <b>genetic algorithm</b> , workspace, <b>parallel manipulator</b> , <b>logic programming</b> , document, grammar, stable model, <b>artificial immune system</b> , logic, word, web site, <b>answer set</b> , global constraint, <b>machining</b> , fitness, belief, evolvability	image, dataset, method, feature, task, <b>convolutional neural network</b> , object, training, <b>deep learning</b> , classification, classifier, <b>deep neural network</b> , <b>neural network</b> , robot, model, video, <b>recurrent neural network</b> , word, <b>segmentation</b> , representation	model, dataset, task, training, <b>transformer</b> , image, <b>deep learning</b> , <b>neural network</b> , prediction, label, <b>federated learning</b> , learning, method, <b>machine learning</b> , <b>language model</b> , explanation, experiment, <b>covid 19</b> , <b>reinforcement learning</b> , feature
PR	problem, <b>algorithm</b> , paper, user, datum, model, word, method, image, information, approach, <b>system</b> , constraint, set, solution, performance, application, document, result, rule	image, algorithm, user, robot, object, network, word, model, dataset, agent, environment, datum, task, video, method, training, <b>language</b> , system, policy, <b>segmentation</b>	image, robot, model, object, algorithm, dataset, task, agent, environment, user, datum, policy, language, graph, <b>reinforcement learning</b> , network, method, video, training, <b>deep learning</b>
FL	paper, problem, <b>algorithm</b> , model, method, approach, <b>system</b> , result, information, set, application, datum, number, user, word, order, performance, case, image, time	image, model, method, paper, problem, approach, datum, algorithm, task, network, performance, dataset, result, user, application, work, feature, system, training, number	image, model, method, paper, task, datum, approach, dataset, performance, problem, training, algorithm, work, result, network, experiment, <b>application</b> , system, feature, <b>deep learning</b>
$KL_{mm}$	problem, paper, <b>algorithm</b> , method, model, <b>system</b> , approach, information, datum, word, set, result, user, application, <b>agent</b> , number, network, performance, <b>language</b> , order	image, model, method, algorithm, datum, paper, task, network, problem, approach, dataset, system, user, feature, performance, training, object, application, result, information	model, method, image, task, datum, dataset, problem, network, approach, paper, training, algorithm, system, performance, feature, object, <b>application</b> , user, <b>deep learning</b> , result
MM	<b>logic program</b> , manipulator, <b>genetic algorithm</b> , workspace, <b>parallel manipulator</b> , <b>logic programming</b> , grammar, <b>stable model</b> , <b>artificial immune system</b> , web site, <b>answer set</b> , <b>global constraint</b> , <b>machining</b> , fitness, <b>evolvability</b> , <b>radial distortion</b> , <b>soft constraint</b> , <b>nonmonotonic reasoning</b> , <b>stable model semantic</b> , <b>belief revision</b>	image, <b>convolutional neural network</b> , <b>recurrent neural network</b> , classifier, <b>deep convolutional neural network</b> , <b>deep network</b> , <b>cnn</b> , <b>computer vision</b> , <b>lstm</b> , <b>deep neural network</b> , <b>bayesian network</b> , <b>rnn</b> , <b>word embedding</b> , <b>svm</b> , <b>segmentation</b> , <b>convolutional network</b> , <b>descriptor</b> , <b>neural machine translation</b> , recognition, sentence	<b>transformer</b> , training, <b>federated learning</b> , <b>language model</b> , <b>covid 19</b> , <b>graph neural network</b> , dataset, explanation, <b>deep learning</b> , <b>pre training</b> , <b>adversarial attack</b> , <b>fine tuning</b> , <b>meta learning</b> , <b>deep learning model</b> , lidar, <b>self attention</b> , point cloud, <b>reinforcement learning</b> , <b>bert</b> , label
$KL_{mm}$	<b>logic program</b> , workspace, <b>genetic algorithm</b> , grammar, manipulator, <b>logic programming</b> , web site, <b>global constraint</b> , <b>artificial immune system</b> , <b>evolvability</b> , <b>parallel manipulator</b> , synonym, <b>stable model</b> , som, <b>belief revision</b> , unification, <b>soft constraint</b> , <b>language resource</b> , fitness, <b>wordnet</b>	image, <b>convolutional neural network</b> , <b>recurrent neural network</b> , classifier, <b>deep network</b> , <b>deep convolutional neural network</b> , <b>bayesian network</b> , <b>word embedding</b> , <b>computer vision</b> , <b>descriptor</b> , <b>svm</b> , <b>crf</b> , <b>lstm</b> , <b>neural machine translation</b> , dictionary, <b>convolutional network</b> , <b>deep neural network</b> , recognition, <b>cnn</b> , <b>segmentation</b>	<b>transformer</b> , training, explanation, <b>language model</b> , <b>covid 19</b> , <b>federated learning</b> , <b>graph neural network</b> , dataset, <b>pre training</b> , lidar, <b>deep learning</b> , <b>adversarial attack</b> , label, <b>meta learning</b> , <b>knowledge distillation</b> , <b>fine tuning</b> , <b>deep learning model</b> , latent space, datum augmentation, target domain

Keywords representative of its corresponding time frame are **highlighted** (annotated by authors).

Table 5: Keyword summaries (top 20 keywords) of three different time frames in AI domain

see  $KL_{mm}$  prefers non-redundant keywords than MM. E.g,  $KL_{mm}$ , instead of selecting *deep neural network* as it already selects keywords like *neural network* and *deep learning*, it selects a different keyword *software* where MM prefers redundant keyword *deep neural network*. Therefore, while the only requirement is to rank keywords based on their distinctiveness for a target domain contrastive with a context domain, MM is more practical to use. On the other hand, if the objective is also selecting diverse representative keywords,  $KL_{mm}$  is preferable. See Appendix D for more qualitative study using newsgroup dataset.

#### 4.2.2 Trending Keywords Selection

As an important application of our problem, we evaluate the performance of proposed approaches for *trending keywords selection* in the AI domain. This study conducts quantitative and qualitative evaluations considering three different time frames: 2000-2009, 2010-2019, and 2020-2021. For this purpose, we compose a corpus representative of each of the specified time frames by collecting abstracts from the Arxiv repository under AI-related categories: cs.AI, cs.CL, cs.CV, cs.IR, cs.LG, cs.NE and cs.RO. The entire dataset under all CS categories is used for the context corpus.

**Quantitative Evaluation.** Since there is no ground truth trending keywords available for the AI domain, it is not easy to quantitatively evaluate the selected ones for a specific time. Instead, we have created three ground truth sets by collecting related keywords from topic areas used in the call for papers (CFP) of an AI conference called AAAI<sup>7</sup> over the three specified time frames. However, the topics that appear in the CFP are not necessarily trending topics, and many topics appear throughout all the time frames. For this, we collect only the changing topics from a time frame to another. Further, to expand the ground truth sets, we also add keywords related to the collected topics. E.g., *word embedding* was a popular keyword in NLP during the 2010s, and one related of this is *word2vec*.

**Evaluation using Ground Truths.** For evaluation, we compute the selected keywords’ ability to cover the ground truth keywords using a *representativeness* measure. Formally, similar to (Kaushal et al., 2019), say  $s_{ij}$  denotes the similarity between two keywords  $t_i$  and  $t_j$ ,  $R(S) = \frac{1}{|\mathcal{G}|} \sum_{t_i \in \mathcal{G}} \max_{t_j \in S} s_{ij}$  is used as the *representativeness* score of selected keyword set  $S$  to represent the ground truth set  $\mathcal{G}$ . For  $s_{ij}$ , we compute the cosine similarity between vector representation of  $t_i$  and  $t_j$ . The vector for

<sup>7</sup><https://www.aaai.org/>



each keyword is the concatenation of two word-vectors; one is word2vec (300d) (Mikolov et al., 2013) learned from corresponding corpus, and the other is the compositional GloVe embedding (Pennington et al., 2014) (element-wise addition of the pre-trained 300d word embeddings). The reason for using pre-trained word vectors is that many keywords in ground truth sets do not appear in the corresponding corpus, and thus vectors cannot be learned from that corpus. Table 3 shows the detailed results over three time frames. We can see that our proposed model  $\mathbf{KL}_{mm}$  outperforms the other methods with large margins followed by MM.

Evaluation using Google Trends. Alongside using ground truths, we also design a quantitative evaluation measure (shown in Table 4) using Google Trends (GT) API<sup>8</sup>. GT<sup>9</sup> awards a score for a term called *interest over time* that expresses the term’s popularity over a specified time range. Since GT does not have data before 2004, we have to use data from 2004 till 2009 for the 2000-2009 time frame. As our three specified time frames are not equal, we first take the average of provided interest scores for each keyword in each time frame to make the score comparable across different time frames. Then, we calculate the probability of each term’s interest over three specified time frames. Finally, the average of computed probability scores of 50 selected terms is calculated for each method. This score represents the average probability of selected terms to be trending in each time frame. From Table 4, we can see that our method  $\mathbf{KL}_{mm}$  achieves the best score over others, followed by comparable results from MM. It indicates that our solutions are more appropriate in finding trending keywords for a specified time frame.

**Qualitative Evaluation.** We qualitatively evaluate the performance of different algorithms by directly comparing their selected keywords in each time frame from Table 5. We can see PR selects keywords that are either CS stop words or the keywords that are not distinctive for a perspective time frame compared to others (similar results by RF, FL,  $\mathbf{KL}_{rf}$ ). Because PR primarily depends on the popularity of a keyword and some keywords always appear frequently in any time frames (e.g., *task*, *dataset*, *model*, etc ). Here, the LO again provides comparable results. E.g., similar to our methods MM and  $\mathbf{KL}_{mm}$ , LO also can select very relevant

trending keywords during the 2020s like *covid 19*. However, while selecting trending keywords, the LO also tends to select many domain-specific stop words overlapped over different time frames (e.g., *method*, *task*, *model*). As discussed before, the reason is that LO does not have the objective of representing the target domain. Therefore, it is not that effective in identifying trending keywords representative for a target domain compared to our models.

## 5 Conclusion

This paper proposes an approach for solving an important but understudied problem of a domain representative keywords selection from candidates contrasting with a context domain. Our approach utilizes a two-component mixture model mechanism followed by a novel subset selection optimization algorithm to tackle the problem. We believe this work will encourage the automated text structuring problem and help a wide range of downstream applications in NLP. For future research direction, we want to focus on adapting the proposed approach in a more challenging task like single document summarization where the scope of information is limited. Besides, our proposed techniques are general and thus can be used in many applications such as information extraction, topic modeling, and concept indexing. Exploration of those applications is an interesting future direction.

## Acknowledgements

We thank the anonymous reviewers for their valuable comments and suggestions. This material is based upon work supported by the National Science Foundation IIS 16-19302 and IIS 16-33755, Zhejiang University ZJU Research 083650, Futurewei Technologies HF2017060011 and 094013, IBM-Illinois Center for Cognitive Computing Systems Research (C3SR)- a research collaboration as part of the IBM Cognitive Horizon Network, grants from eBay and Microsoft Azure, UIUC OVCR CCIL Planning Grant 434S34, UIUC CSBS Small Grant 434C8U, and UIUC New Frontiers Initiative. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect the views of the funding agencies.

<sup>8</sup><https://github.com/GeneralMills/pytrends>

<sup>9</sup><https://trends.google.com>

## References

- Rabah Alzaidy, Cornelia Caragea, and C Lee Giles. 2019. Bi-lstm-crf sequence labeling for keyphrase extraction from scholarly documents. In *The world wide web conference*, pages 2551–2557.
- 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*.
- Adam L. Berger and John D. Lafferty. 1999. **Information retrieval as statistical translation**. In *SIGIR '99: Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, August 15-19, 1999, Berkeley, CA, USA*, pages 222–229. ACM.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022.
- Gerlof Bouma. 2009. Normalized (pointwise) mutual information in collocation extraction. *Proceedings of GSCL*, 30:31–40.
- Peter F Brown, Stephen A Della Pietra, Vincent J Della Pietra, and Robert L Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational linguistics*, 19(2):263–311.
- Sheng Chen, Akshay Soni, Aasish Pappu, and Yashar Mehdad. 2017. Doctag2vec: An embedding based multi-label learning approach for document tagging. *arXiv preprint arXiv:1707.04596*.
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. 1977. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22.
- Uriel Feige. 1998. A threshold of  $\ln n$  for approximating set cover. *Journal of the ACM (JACM)*, 45(4):634–652.
- Abram Handler, Matthew Denny, Hanna Wallach, and Brendan O’Connor. 2016. Bag of what? simple noun phrase extraction for text analysis. In *Proceedings of the First Workshop on NLP and Computational Social Science*, pages 114–124.
- Anna Häty, Michael Dorna, and Sabine Schulte im Walde. 2017. Evaluating the reliability and interaction of recursively used feature classes for terminology extraction. In *Proceedings of the student research workshop at the 15th conference of the European chapter of the association for computational linguistics*, pages 113–121.
- Thomas Hofmann. 2001. Unsupervised learning by probabilistic latent semantic analysis. *Machine learning*, 42(1):177–196.
- Jie Huang, Kevin Chen-Chuan Chang, Jinjun Xiong, and Wen-mei Hwu. 2021. Measuring fine-grained domain relevance of terms: A hierarchical core-fringe approach. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*.
- Jack Hughes, Seth Aycock, Andrew Caines, Paula Buttery, and Alice Hutchings. 2020. Detecting trending terms in cybersecurity forum discussions. In *Proceedings of the Sixth Workshop on Noisy User-generated Text (W-NUT 2020)*, pages 107–115.
- Maryam Karimzadehgan and ChengXiang Zhai. 2010. Estimation of statistical translation models based on mutual information for ad hoc information retrieval. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*, pages 323–330.
- Vishal Kaushal, Rishabh Iyer, Suraj Kothawade, Rohan Mahadev, Khoshrav Doctor, and Ganesh Ramakrishnan. 2019. Learning from less data: A unified data subset selection and active learning framework for computer vision. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1289–1299. IEEE.
- Solomon Kullback and Richard A Leibler. 1951. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86.
- Jialu Liu, Jingbo Shang, Chi Wang, Xiang Ren, and Jiawei Han. 2015. Mining quality phrases from massive text corpora. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pages 1729–1744.
- Weiming Lu, Yangfan Zhou, Jiale Yu, and Chenhao Jia. 2019. Concept extraction and prerequisite relation learning from educational data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9678–9685.
- Rui Meng, Sanqiang Zhao, Shuguang Han, Daqing He, Peter Brusilovsky, and Yu Chi. 2017. Deep keyphrase generation. *arXiv preprint arXiv:1704.06879*.
- Rada Mihalcea and Paul Tarau. 2004. Textrank: Bringing order into text. In *Proceedings of the 2004 conference on empirical methods in natural language processing*, pages 404–411.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Michel Minoux. 1978. Accelerated greedy algorithms for maximizing submodular set functions. In *Optimization techniques*, pages 234–243. Springer.
- Pitu B Mirchandani and Richard L Francis. 1990. *Discrete location theory*.

- Burt L Monroe, Michael P Colaresi, and Kevin M Quinn. 2008. Fightin' words: Lexical feature selection and evaluation for identifying the content of political conflict. *Political Analysis*, 16(4):372–403.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Sunita Sarawagi, Soumen Chakrabarti, and Shantanu Godbole. 2003. Cross-training: Learning probabilistic mappings between topics. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 177–186.
- Jingbo Shang, Jialu Liu, Meng Jiang, Xiang Ren, Clare R Voss, and Jiawei Han. 2018. Automated phrase mining from massive text corpora. *IEEE Transactions on Knowledge and Data Engineering*, 30(10):1825–1837.
- Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. 2008. Arnetminer: extraction and mining of academic social networks. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 990–998.
- Peter D Turney. 2000. Learning algorithms for keyphrase extraction. *Information retrieval*, 2(4):303–336.
- Li Wang, Wei Zhu, Sihang Jiang, Sheng Zhang, Keqiang Wang, Yuan Ni, Guotong Xie, and Yanghua Xiao. 2020. Mining infrequent high-quality phrases from domain-specific corpora. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 1535–1544.
- Ian H. Witten, Gordon W. Paynter, Eibe Frank, Carl Gutwin, and Craig G. Nevill-Manning. 1999. **KEA: practical automatic keyphrase extraction**. In *Proceedings of the Fourth ACM conference on Digital Libraries, August 11-14, 1999, Berkeley, CA, USA*, pages 254–255. ACM.
- Chengxiang Zhai and John Lafferty. 2001. Model-based feedback in the language modeling approach to information retrieval. In *Proceedings of the tenth international conference on Information and knowledge management*, pages 403–410.
- ChengXiang Zhai, Atulya Velivelli, and Bei Yu. 2004. A cross-collection mixture model for comparative text mining. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 743–748.
- Chao Zhang, Fangbo Tao, Xiusi Chen, Jiaming Shen, Meng Jiang, Brian Sadler, Michelle Vanni, and Jiawei Han. 2018. Taxogen: Constructing topical concept taxonomy by adaptive term embedding and clustering. *Proc. KDDI*.

$k$	RF	LO	PR	FL	$KL_{rf}$	MM	$KL_{mm}$
Candidate Keywords from Springer							
10	2.0691	2.1157	1.1273	1.0852	1.0998	2.1401	<b>2.1432</b>
20	2.1745	2.2391	2.2557	2.1745	2.1805	2.2915	<b>3.3260</b>
30	2.2302	3.3938	3.4220	2.2261	2.2556	3.4259	<b>3.4262</b>
40	2.2873	3.5040	3.5053	2.2847	2.3223	5.5379	<b>5.5404</b>
50	2.3591	5.6136	3.6136	2.3591	2.3689	5.6346	<b>5.6576</b>
100	3.7701	7.0986	6.0387	3.7601	3.7846	<b>8.1300</b>	7.1193
200	6.4205	12.0418	8.8759	5.3433	6.3992	<b>12.1043</b>	12.0591
500	13.3390	19.4237	13.9664	13.2692	12.2365	<b>19.4987</b>	19.4857
Candidate Keywords from Aminer							
10	1.1020	2.1205	1.1101	1.1073	1.1109	2.1205	<b>2.1645</b>
20	2.1782	2.3050	1.2272	2.1699	2.1956	3.3048	<b>3.3079</b>
30	2.2433	3.4320	2.3747	2.2508	2.2795	3.4394	<b>4.4522</b>
40	2.3159	3.5312	3.5345	2.3012	2.3944	5.5693	<b>5.5700</b>
50	2.4617	5.6519	4.6396	3.4657	2.4675	5.6463	<b>5.6829</b>
100	3.8432	8.1523	6.1264	3.8038	3.8429	<b>9.1788</b>	8.1679
200	6.5778	12.1533	8.9075	6.5320	6.5431	<b>12.2007</b>	12.1013
500	13.668	19.5232	15.1445	13.4622	12.4942	19.5580	<b>20.5255</b>

Table 6: Results of selected summary keywords’ correspondence with arxiv category keywords

## A Proof of Theorem 1

To prove this, we need to first show that  $\phi(\cdot)$  is *submodular* and *monotone*. As, we are concerned on minimizing  $\phi(\cdot)$ , it is equivalent to maximizing  $F(\cdot) = -\phi(\cdot)$ . Hence, it is sufficient to prove that  $F(\cdot)$  is submodular and monotone. Let,  $X \subseteq Y \subseteq V$  and  $v \in V$ , then we get

$$\begin{aligned}
& F(X \cup \{v\}) - F(\{X\}) = \\
& \sum_{t_i \in V} p(t_i) \log \frac{\sum_{t_j \in XU\{v\}} p(t_i|t_j)p(t_j)}{p(t_i)} \\
& - \sum_{t_i \in V} p(t_i) \log \frac{\sum_{t_j \in X} p(t_i|t_j)p(t_j)}{p(t_i)} \\
& = \sum_{t_i \in V} p(t_i) \log \frac{\sum_{t_j \in XU\{v\}} p(t_i|t_j)p(t_j)}{\sum_{t_j \in X} p(t_i|t_j)p(t_j)} \\
& = \sum_{t_i \in V} p(t_i) \log \frac{\sum_{t_j \in X} p(t_i|t_j)p(t_j) + p(t_i|v)p(v)}{\sum_{t_j \in X} p(t_i|t_j)p(t_j)} \\
& = \sum_{t_i \in V} p(t_i) \log \left( 1 + \frac{p(t_i|v)p(v)}{\sum_{t_j \in X} p(t_i|t_j)p(t_j)} \right).
\end{aligned}$$

Similarly,  $F(Y \cup \{v\}) - F(\{Y\}) =$

$$\sum_{t_i \in V} p(t_i) \log \left( 1 + \frac{p(t_i|v)p(v)}{\sum_{t_j \in Y} p(t_i|t_j)p(t_j)} \right).$$

As  $X \subseteq Y$ , then,

$$\sum_{t_j \in X} p(t_i|t_j)p(t_j) \leq \sum_{t_j \in Y} p(t_i|t_j)p(t_j).$$

Therefore,  $F(X \cup \{v\}) - F(\{X\}) \geq F(Y \cup \{v\}) - F(\{Y\})$  which proves that  $F(\cdot)$  is *submodular*. Moreover, we can show that  $F(Y) - F(X) =$

Selected 20 keywords in CS	
RF	model, method, paper, problem, approach, <b>algorithm</b> , datum, <b>network</b> , <b>system</b> , <b>performance</b> , task, result, <b>image</b> , number, user, application, time, <b>dataset</b> , <b>graph</b> , work
LO	task, <b>algorithm</b> , user, <b>network</b> , <b>performance</b> , <b>dataset</b> , <b>image</b> , problem, <b>training</b> , approach, <b>deep learning</b> , method, <b>node</b> , <b>agent</b> , language, <b>neural network</b> , paper, video, challenge, <b>architecture</b>
PR	<b>image</b> , <b>graph</b> , <b>dataset</b> , user, method, model, <b>network</b> , task, <b>algorithm</b> , datum, problem, <b>system</b> , training, channel, node, performance, object, agent, deep learning, language
FL	<b>image</b> , model, paper, problem, method, network, datum, approach, <b>algorithm</b> , <b>system</b> , user, performance, result, <b>graph</b> , task, application, number, work, time, dataset
$KL_{rf}$	model, method, problem, <b>system</b> , <b>network</b> , datum, <b>algorithm</b> , paper, <b>image</b> , task, user, approach, performance, graph, application, dataset, time, result, number, information
MM	task, <b>algorithm</b> , user, <b>network</b> , <b>performance</b> , <b>dataset</b> , image, <b>training</b> , <b>deep learning</b> , <b>node</b> , <b>agent</b> , language, <b>neural network</b> , video, <b>architecture</b> , challenge, <b>robot</b> , real world, <b>attack</b> , <b>learning</b>
$KL_{mm}$	task, user, <b>algorithm</b> , <b>dataset</b> , <b>network</b> , <b>performance</b> , <b>image</b> , <b>training</b> , <b>agent</b> , language, <b>deep learning</b> , <b>attack</b> , <b>robot</b> , video, challenge, <b>node</b> , <b>neural network</b> , <b>query</b> , <b>code</b> , <b>machine learning</b>

Keywords distinctive to the CS domain are **highlighted** (annotated by authors).

Table 7: Summary keywords selected from Springer candidate keywords in CS domain

Selected 20 keywords in CS	
RF	model, method, <b>network</b> , <b>algorithm</b> , <b>system</b> , datum, problem, user, <b>image</b> , time, <b>graph</b> , application, <b>performance</b> , state, feature, <b>dataset</b> , number, art, work, information
LO	<b>network</b> , user, <b>algorithm</b> , <b>dataset</b> , art, <b>performance</b> , training, <b>image</b> , task, <b>deep learning</b> , <b>node</b> , <b>learning</b> , <b>agent</b> , <b>attack</b> , <b>neural network</b> , language, video, problem, <b>robot</b> , <b>graph</b>
PR	<b>image</b> , art, <b>graph</b> , <b>dataset</b> , state, user, model, <b>network</b> , method, <b>algorithm</b> , datum, vertex, training, channel, <b>system</b> , feature, <b>node</b> , <b>deep learning</b> , experiment, object
FL	art, model, method, <b>network</b> , <b>algorithm</b> , <b>image</b> , datum, <b>system</b> , problem, user, <b>graph</b> , application, <b>performance</b> , time, work, number, <b>dataset</b> , feature, order, experiment
$KL_{rf}$	model, method, <b>algorithm</b> , <b>network</b> , <b>system</b> , datum, <b>image</b> , user, problem, <b>graph</b> , <b>dataset</b> , application, <b>performance</b> , time, <b>feature</b> , <b>agent</b> , art, information, number, <b>training</b>
MM	<b>network</b> , user, <b>algorithm</b> , <b>dataset</b> , art, <b>performance</b> , <b>training</b> , <b>image</b> , task, <b>deep learning</b> , <b>node</b> , <b>learning</b> , <b>agent</b> , <b>neural network</b> , <b>attack</b> , language, video, <b>robot</b> , <b>architecture</b> , <b>machine learning</b>
$KL_{mm}$	<b>dataset</b> , user, <b>network</b> , <b>algorithm</b> , <b>training</b> , <b>image</b> , <b>agent</b> , task, <b>performance</b> , <b>attack</b> , art, language, <b>deep learning</b> , <b>robot</b> , <b>learning</b> , video, <b>node</b> , <b>machine learning</b> , <b>code</b> , <b>neural network</b>

Keywords distinctive to the CS domain are **highlighted** (annotated by authors).

Table 8: Summary keywords selected from Aminer candidate keywords in CS domain

Religion	Recreation	Science	Politics
talk.religion.misc	rec.autos	sci.crypt	talk.politics.misc
alt.atheism	rec.motorcycles	sci.electronics	talk.politics.guns
soc.religion.christian	rec.sport.baseball	sci.med	talk.politics.mideast
	rec.sport.hockey	sci.space	

Table 9: Subtopics in each of the four topics in the newsgroup dataset

$\sum_{t_i \in V} p(t_i) \log \frac{\sum_{t_j \in Y} p(t_i|t_j)p(t_j)}{\sum_{t_j \in X} p(t_i|t_j)p(t_j)} \geq 1$ . Hence,  $F(Y) \geq F(X)$  for  $X \subseteq Y \subseteq V$  which proves that  $F(\cdot)$  is monotone. Therefore, it proves that minimizing  $\phi(\cdot)$  using simple greedy algorithm guarantees an approximation of  $1 - \frac{1}{e}$ .

	Religion	Recreation	Science	Politics
RF	thing, <b>church</b> , life, word, man, <b>religion</b> , <b>bible</b> , <b>faith</b> , question, <b>belief</b> , book, point, law, evidence, <b>sin</b> , reason, world, truth, child, <b>god</b>	<b>game</b> , <b>car</b> , <b>team</b> , <b>player</b> , <b>bike</b> , <b>season</b> , <b>point</b> , <b>hockey</b> , problem, lot, <b>goal</b> , <b>baseball</b> , guy, <b>engine</b> , power, number, year, line, question, <b>run</b>	key, information, thing, government, <b>space</b> , <b>encryption</b> , datum, <b>clipper</b> , <b>chip</b> , case, number, phone, <b>bit</b> , <b>privacy</b> , drug, earth, power, <b>security</b> , <b>program</b> , disease	<b>government</b> , <b>gun</b> , child, state, law, country, man, <b>president</b> , case, <b>war</b> , group, fact, <b>firearm</b> , number, <b>crime</b> , question, <b>weapon</b> , world, history, population
LO	<b>church</b> , <b>bible</b> , <b>faith</b> , <b>religion</b> , <b>belief</b> , <b>sin</b> , <b>god</b> , scripture, life, word, <b>atheist</b> , truth, <b>atheism</b> , <b>homosexuality</b> , love, man, evidence, son, <b>morality</b> , book	<b>game</b> , <b>team</b> , <b>car</b> , <b>player</b> , <b>bike</b> , <b>season</b> , <b>hockey</b> , <b>baseball</b> , <b>playoff</b> , <b>engine</b> , <b>goal</b> , <b>pitcher</b> , <b>tire</b> , <b>run</b> , pen, <b>league</b> , <b>puck</b> , <b>motorcycle</b> , dog, <b>clutch</b>	key, <b>encryption</b> , <b>space</b> , <b>clipper</b> , <b>privacy</b> , <b>satellite</b> , <b>mission</b> , <b>disease</b> , <b>shuttle</b> , phone, <b>orbit</b> , <b>escrow</b> , <b>moon</b> , <b>cancer</b> , <b>algorithm</b> , <b>spacecraft</b> , <b>security</b> , launch, <b>vitamin</b> , health	<b>gun</b> , <b>government</b> , <b>firearm</b> , <b>president</b> , country, <b>weapon</b> , <b>crime</b> , village, <b>soldier</b> , <b>genocide</b> , <b>war</b> , population, state, child, police, <b>turk</b> , <b>massacre</b> , <b>handgun</b> , compound, new york
PR	man, word, thing, life, world, history, <b>church</b> , book, question, <b>bible</b> , <b>faith</b> , point, truth, reason, matter, law, year, <b>religion</b> , earth, mind	<b>game</b> , <b>team</b> , <b>player</b> , <b>car</b> , <b>season</b> , <b>goal</b> , <b>point</b> , <b>hockey</b> , shot, year, power, number, <b>engine</b> , <b>bike</b> , <b>win</b> , <b>league</b> , speed, line, <b>run</b> , end	information, datum, year, study, number, united states, <b>space</b> , nature, <b>security</b> , mail, government, thing, <b>encryption</b> , case, archive, key, life, book, law, <b>science</b>	<b>government</b> , man, group, <b>war</b> , world, village, child, fact, year, history, life, state, house, end, woman, home, <b>power</b> , arm, law, population
FL	man, thing, <b>church</b> , life, word, <b>religion</b> , book, question, history, point, <b>bible</b> , law, evidence, reason, <b>sin</b> , world, <b>belief</b> , child, <b>faith</b> , case	<b>game</b> , <b>car</b> , <b>team</b> , <b>player</b> , <b>bike</b> , <b>season</b> , pit, problem, <b>hockey</b> , lot, power, <b>point</b> , <b>baseball</b> , <b>engine</b> , <b>goal</b> , <b>run</b> , question, guy, standing, <b>speed</b>	information, <b>space</b> , key, study, thing, government, datum, year, case, number, patient, <b>software</b> , power, book, archive, food, <b>clipper</b> , <b>mission</b> , hicnet medical newsletter page, <b>encryption</b>	<b>government</b> , man, village, <b>gun</b> , <b>president</b> , <b>sumgait</b> , history, state, case, child, <b>law</b> , world, country, population, fact, <b>war</b> , los angeles, number, group, year
KL <sub>rf</sub>	thing, life, <b>church</b> , word, <b>belief</b> , man, question, <b>religion</b> , <b>bible</b> , <b>faith</b> , book, law, evidence, point, <b>sin</b> , world, reason, child, truth, <b>god</b>	<b>game</b> , <b>car</b> , <b>team</b> , <b>bike</b> , <b>player</b> , <b>point</b> , problem, <b>season</b> , lot, <b>hockey</b> , <b>engine</b> , guy, power, <b>baseball</b> , question, goal, number, list, road, year	key, information, <b>space</b> , thing, government, <b>encryption</b> , datum, case, <b>clipper</b> , number, <b>chip</b> , <b>disease</b> , power, phone, <b>earth</b> , drug, <b>program</b> , <b>bit</b> , book, <b>privacy</b>	<b>government</b> , <b>gun</b> , child, state, country, man, law, president, <b>war</b> , case, group, fact, question, <b>crime</b> , population, number, world, <b>firearm</b> , history, woman
MM	<b>bible</b> , <b>church</b> , <b>faith</b> , <b>sin</b> , scripture, <b>atheism</b> , <b>god</b> , <b>gospel</b> , christianity, <b>prophecy</b> , <b>mcconkie</b> , <b>jesus christ</b> , <b>prophet</b> , new testament, <b>atheist</b> , <b>disciple</b> , holy spirit, <b>theist</b> , christian, <b>homosexuality</b>	team, pen, <b>player</b> , <b>bike</b> , <b>hockey</b> , <b>season</b> , second period, <b>puck</b> , <b>playoff</b> , first period, <b>ranger</b> , schedule, <b>pitcher</b> , <b>baseball</b> , <b>nhl</b> , <b>cub</b> , <b>tire</b> , <b>injury</b> , <b>league</b> , respect	<b>encryption</b> , <b>clipper</b> , <b>privacy</b> , <b>satellite</b> , <b>shuttle</b> , <b>orbit</b> , <b>vitamin</b> , <b>infection</b> , <b>escrow</b> , <b>moon</b> , <b>pgp</b> , <b>mission</b> , <b>spacecraft</b> , <b>cryptography</b> , <b>cancer</b> , <b>circuit</b> , <b>astronaut</b> , <b>asteroid</b> , <b>cipher</b> , <b>telescope</b>	<b>gun</b> , <b>firearm</b> , <b>soldier</b> , village, <b>genocide</b> , <b>bayonet</b> , <b>turk</b> , <b>handgun</b> , <b>massacre</b> , new york, <b>tartar</b> , <b>homicide</b> , <b>civilian</b> , <b>weapon</b> , <b>human right</b> , <b>gun control</b> , <b>bullet</b> , <b>troop</b> , <b>ottoman</b> , <b>sumgait</b>
KL <sub>mm</sub>	<b>bible</b> , <b>church</b> , <b>faith</b> , <b>sin</b> , scripture, <b>atheism</b> , <b>god</b> , <b>gospel</b> , <b>jesus christ</b> , <b>prophecy</b> , christianity, <b>new testament</b> , <b>mcconkie</b> , <b>prophet</b> , <b>holy spirit</b> , <b>morality</b> , <b>theist</b> , <b>disciple</b> , <b>homosexuality</b> , <b>atheist</b>	team, pen, <b>bike</b> , <b>player</b> , <b>season</b> , <b>hockey</b> , second period, <b>playoff</b> , <b>pitcher</b> , <b>puck</b> , schedule, <b>ranger</b> , <b>cub</b> , <b>baseball</b> , <b>tire</b> , <b>clutch</b> , first period, favor, respect, <b>nhl</b>	<b>encryption</b> , <b>satellite</b> , <b>clipper</b> , <b>vitamin</b> , <b>shuttle</b> , <b>infection</b> , <b>orbit</b> , <b>moon</b> , <b>cancer</b> , <b>privacy</b> , <b>spacecraft</b> , <b>circuit</b> , <b>mission</b> , <b>pgp</b> , <b>escrow</b> , <b>allergy</b> , <b>yeast</b> , <b>cryptography</b> , <b>diet</b> , <b>solar sail</b>	<b>gun</b> , village, <b>firearm</b> , <b>soldier</b> , <b>genocide</b> , <b>turk</b> , new york, <b>massacre</b> , <b>human right</b> , <b>handgun</b> , <b>bayonet</b> , <b>civilian</b> , <b>croat</b> , <b>tartar</b> , <b>weapon</b> , <b>gun control</b> , <b>troop</b> , <b>homicide</b> , well regulated, <b>sumgait</b>

Keywords distinctive to the subtopics in a respected topic are **highlighted** (annotated by authors).

Table 10: Summary keywords selected by different algorithms on four topics from newsgroups dataset

## **B Additional Quantitative Results on CS Domain**

Results are shown in Table 6.

## **C Additional Qualitative Results on CS Domain**

Results are shown in Tables 7 and 8.

## **D Evaluation Using Newsgroup Dataset**

We use newsgroup dataset covering four known topics named *Religion*, *Recreation*, *Science* and *Politics*. In this study, we split the whole dataset into these four topic groups represented by their corpus and use the whole newsgroup dataset as our background corpus. Table 9 shows the subtopics for each of the four topics. For each topic, we show the selected top 20 keywords using different algorithms in Table 10. This study aims to evaluate the capability of the proposed models to select distinctive keywords for each topic compared to the baselines. We can see that almost all the keywords selected by our methods MM and  $\mathbf{KL}_{mm}$  are distinctive for each topic relating closely with respected subtopics shown in Table 9 and do not overlap with other topics. Similarly, as previously, the results from LO come close to ours with some anomalies. For instance, our methods select informative keywords like *jesus christ*, *holy spirit* and *new testament* for *religion* topic rather than non-distinctive keywords like *word*, *man* and *son*.