# KS-GNN: Keywords Search over Incomplete Graphs via Graphs Neural Network

Anonymous Author(s) Affiliation Address email

#### Abstract

Keyword search is a fundamental task to retrieve information that is the most 1 2 relevant to the query keywords. Keyword search over graphs aims to find subtrees 3 or subgraphs containing all query keywords ranked according to some criteria. Existing studies all assume that the graphs have complete information. However, 4 5 real-world graphs may usually contain some missing information (such as edges or keywords), thus making the problem much more challenging. To solve the problem 6 of keyword search over incomplete graphs, we propose a novel model named 7 KS-GNN based on the graph neural network and the auto-encoder. By considering 8 9 the latent relationships and the frequency of different keywords, the proposed KS-GNN aims to alleviate the effect of missing information and is able to learn 10 low-dimensional representative node embeddings that preserve both graph structure 11 and keyword features. Our model can effectively answer keyword search queries 12 with linear time complexity over incomplete graphs. The experiments on four 13 real-world datasets show that our model consistently achieves better performance 14 than state-of-the-art baseline methods in graphs having missing information. 15

### 16 **1** Introduction

Keyword search is an important research topic 17 which allows users to provide query keywords 18 and returns the most relevant results. The key-19 word search over graph data [1] usually retrieves 20 top-k subtrees or subgraphs which contain all 21 the query keywords ranked according to some 22 criteria. For example, He et al. [2] propose a 23 general scoring function considering both graph 24 structure and content, and they aim to find top-k25 nodes where each node can reach all query key-26 words, and the sum of its shortest path distances 27 to these keywords is as small as possible. This 28 ranking method is commonly used in later graph 29 keyword search works [3, 4]. 30



Figure 1: Example of keyword search on incomplete graphs

All existing studies assume that the graph data is complete and has no missing information. However, real-world graphs may usually have some missing edges [5] and missing attributes on some nodes [6]. This renders previous graph keyword search methods fail in finding exact answers when dealing with such incomplete graphs. Figure 1 shows an example keyword search query over graphs. Given  $q=\{c, e, f\}$ , on the left graph G with no missing information, the best node is  $v_4$ , since it contains keywords c and f, it can reach  $v_5$  containing e, and its sum of the shortest path distances to all

Submitted to 35th Conference on Neural Information Processing Systems (NeurIPS 2021). Do not distribute.

query keywords is the smallest which is 1 (the shortest distances of  $v_4$  to c, e, and f are 0, 1, and 0). However, on the right graph G' which has a missing edge and a node with missing attributes, the

result becomes  $v_1$ ,  $v_2$  or  $v_6$ , and the subtree consists of  $\{v_1, v_2, v_6\}$ , with a total distance of 2.

To handle the missing information, one simple idea is to first utilize some state-of-the-art graph 40 completion models (such as SAT [6]) to predict the missing information and then apply the existing 41 algorithms (such as BLINKS [2]) to find the answers from the graph with predicted keywords and 42 edges. However, such a completed graph contains many noises and errors comparing with the original 43 graph, and thus this method has poor performance as shown in our experimental study. To capture 44 the latent information of the incomplete graphs, we propose to utilize the graph neural network 45 (GNN) for graph keyword search. GNN has been widely applied in tasks such as link prediction, 46 node classification, and node clustering [7, 6, 8], but existing models cannot be directly applied to 47 keyword search since they usually embed all the features (keywords) of a node into a single vector 48 and they cannot obtain the representation for the individual query keywords. 49

We firstly design two naive approaches based on GNN and dimension reduction. To achieve better 50 performance, we propose a novel auto-encoder and GNN-based model using the message passing 51 mechanism, called KS-GNN. The model mainly consists of three components: an encoder that 52 transforms the original keyword information to low-dimensional embedding vectors; a decoder 53 that aims to reconstruct the high-dimensional representation of keywords from the embedding; a 54 message passing-based aggregation mechanism that preserves the shortest path information between 55 keywords and the target node. Different from the existing graph keyword search works, we propose to 56 leverage GNN to obtain representative node embedding that contains the keyword information, taking 57 the latent graph structure, keyword distribution, and keyword frequency information into account. 58 Meantime, the proposed KS-GNN is able to encode the input query as a low-dimensional vector by 59 its learned powerful encoder, and the results are obtained by computing the similarity between the 60 query embedding and node embeddings. This also speeds up the query processing to the cost of linear 61 time complexity. 62

<sup>63</sup> The main contributions of our approach are as follows:

- To our best knowledge, this is the first work on keyword search in graphs with missing information.
- We propose an auto-encoder and GNN-based model KS-GNN to solve the problem effectively without having to know the complete information of the input graph.
- The experimental results on several real-world datasets show that our proposed model consistently outperforms several baseline methods.

## 70 2 Related Work

**Keyword Search in Graphs.** Keyword search over graph data aims to find the top-k subtrees or 71 subgraphs according to some ranking criteria. The conventional methods design algorithms assuming 72 that the graphs have complete information. For example, DBXplorer [9] proposes to utilize the 73 number of the answer's edges as the scoring function. BANKS [10] model tuples as nodes in a graph 74 and then performs keyword search using proximity-based ranking. He et al. [2] proposes a general 75 ranking function considering both graph structure and content. BLINKS also builds an efficient 76 bi-level index structure to improve efficiency. Kargar and An, motivated by the Steiner tree problem, 77 use the total edge weight in ranking answers [11]. There also exists studies on keyword search in 78 temporal graphs [12], uncertain graphs [13], knowledge graphs [14], RDF graphs [15], etc. However, 79 real-world graphs may usually be incomplete. As all the state-of-art keyword search methods retrieve 80 the exact answer, the missing information (keywords or edges) imposes a significant effect on the 81 query results. To address this issue, we propose a graph representation learning-based solution to 82 solve the top-k keyword search problem on incomplete graphs. 83

Graph Neural Networks. As a powerful branch of graph representation learning methods, the
graph neural network has been widely used in recent years due to its excellent performance. The
models in the early years are usually based on the so-called graph convolutional network (GCN) [16,
17, 18], which is based on the Fourier transform theory of graphs developed by Shuman et al. [19].
However, research in recent years has shown that the GCN-based methods can be represented by

the message passing mechanism, which is more consistent with the experimental results [20]. The 89 Graph Attention Network (GAT) [21] is one of the representatives of graph neural networks based on 90 the message passing mechanism. GAT introduces the attention mechanism to calculate the attention 91 coefficient between nodes and then uses it to assign different weights to neighbors' information. 92 Based on the auto-encoder and GCN, Graph auto-encoder (GAE) [22] is proposed to reconstruct 93 the adjacency matrix. Moreover, there are some GNN-based works that aim to predict and impute 94 missing data to a data matrix [23, 24, 6]. However, all the methods mentioned above cannot directly 95 handle the graph keyword search problem. To our best knowledge, this is the first work that leverages 96 GNN to process keyword search on incomplete graphs. 97

# 98 **3** Problem Statement

A graph keyword search query  $q = (w_{q_1}, w_{q_2}, ..., w_{q_m})$  contains a set of query keywords, and it searches relevant results from a graph<sup>1</sup>  $G = (\mathcal{V}, \mathcal{E}, \mathcal{W})$ , where each node  $v \in \mathcal{V}$ , each edge  $e \in \mathcal{E}$ , 99 100 and each keyword  $w \in \mathcal{W}$ . For each node v, it is associated with a set of keywords  $\{w_1^v, w_2^v, ..., w_n^v\}$ . 101 In this work, we study the keyword search problem over an incomplete graph. To alleviate the effect 102 of the missing information to keyword search over incomplete graphs, we assume that in the original 103 graph the query results are obtained by applying the BLINKS scheme [2] (a commonly used graph 104 keyword search method). Given a query q, let s(v,q) denote the score of the node v. According 105 to [2, 3, 25, 26, 27],  $s(v,q) = \sum_{i=1}^{m} dist_{min}(v, w_{q_i})$ , where  $dist_{min}(v, w_{q_i})$  computes the shortest path distance from node v to a node containing  $w_{q_i}$ . BLINKS and to ind top-k nodes where each 106 107 node can reach all query keywords in the graph, and the scores of the k nodes measured by s(v, q) are 108 the smallest. E.g., in Figure 1(a),  $s(v_4, q) = dist_{min}(v_4, c) + dist_{min}(v_4, e) + dist_{min}(v_4, f) = 1$ . 109

**Problem Definition.** Given an incomplete graph  $G' = (\mathcal{V}, \mathcal{E}', \mathcal{W}', r_w, r_e)$ , where  $\mathcal{E}' \subseteq \mathcal{E}, \mathcal{W}' \subseteq \mathcal{W}$ , and the proportions of nodes with missing keywords and of missing edges in G are denoted by  $r_w$ and  $r_e$ , respectively. Given a query q, the incomplete graph top-k keyword search problem aims to find a set  $S = (v_1, v_2, ..., v_k)$  of k nodes from G' such that for any vertex  $v' \notin S$ ,  $s(v', q) \ge$  $max(\{s(v_i, q) | v_i \in S\})$ .

## 115 4 Proposed Methods

We propose to solve the keyword search prob-116 lem with an unsupervised graph representation 117 learning method, since the representative low-118 dimensional node embeddings can capture the 119 latent information of the input incomplete graph 120 and thus can help recovering the missing infor-121 mation. In addition, low dimensional node em-122 beddings can speed up the query processing by 123 comparing the node embedding with the gen-124 erated query embedding at the cost of linear 125 complexity. In this section, we first propose two 126 naive methods that are based on GNN and di-127 mensionality reduction methods and then we 128 introduce our proposed KS-GNN in details. 129



Figure 2: An illustration of the message passing and aggregation of Conv-OH, where  $v_3$  is a target node and it aggregates keyword information from its neighbors.

#### 130 4.1 Naive Methods

Conv-OH. Graph convolutional layer has been widely used in GNNs, which enables GNN models
 to gather information from neighbor nodes. Our first naive method Conv-OH utilizes the graph
 convolutional layer and takes the one-hot encoding of keywords as the input feature for nodes.

Specifically, with  $|\mathcal{V}| = N$  and  $|\mathcal{W}| = M$ , the one-hot encoding of node v is denoted by  $\mathbf{x} = \{0, 1\}^M$ with  $h_{v,j} = 1$  if  $w_j$  is a keyword associated with v and 0 otherwise. Therefore, the input feature matrix is denoted by  $\mathbf{X} \in \{0, 1\}^{N \times M}$ . Let  $\mathbf{H}^l$  denote the output node embedding of the *l*-th layer,

<sup>&</sup>lt;sup>1</sup>For ease of presentation, we focus on the undirected graphs, and it is easy to extend the proposed method in directed graphs by passing messages along the edges.

137 and we have:

$$\mathbf{h}_{v}^{l+1} = Aggregate\Big(\{f(\mathbf{h}_{u}^{l}), \forall u \in \mathcal{N}(v)\} \cup \{\mathbf{h}_{v}^{l}\}\Big),\tag{1}$$

where  $\mathbf{H}^0 = \mathbf{X}$ ,  $f(\cdot)$  denotes a transform function and  $\mathcal{N}(v)$  denotes the neighbors of node v. Using the combined distance as the scoring function (described in Section 3), Eq. (1) can be written as:

$$\mathbf{h}_{v}^{l+1} = \Omega\Big(\{sgn(\mathbf{h}_{u}^{l}) \circ (\mathbf{h}_{u}^{l}+1), \forall u \in \mathcal{N}(v)\} \cup \{\mathbf{h}_{v}^{l}\}\Big),\tag{2}$$

where  $\Omega(\cdot)$  denotes the element-wise minimum function which ignores zeros, and  $sgn(\mathbf{h})$  denotes the sign function that returns a vector with the signs of the corresponding elements of  $\mathbf{h}$  (the sign of an element is 1 if the element is positive and 0 otherwise). For instance,  $\Omega(\{[0, 0, 1, 1], [2, 0, 0, 2]\}) =$ [2, 0, 1, 1], and sgn([0, 2, 0, 3]) = [0, 1, 0, 1].

With Eq. (2), the output of Conv-OH is an  $N \times M$  matrix, denoted by **Z**. Note that there is no dimensionality reduction in Conv-OH and thus this method consumes huge space. The node embedding  $\mathbf{h}_{v}^{l}$ of v also represents the shortest path distances between the keywords and v. Specifically, if  $h_{v,i}^{l} > 0$ , it means that the shortest path distance between v and  $w_{i}$  is  $h_{v,i}^{l} - 1$ , and v cannot reach  $w_{i}$  within l hops if  $h_{v,i}^{l} = 0$ . Hence, Conv-OH is able to return the same answer as does BLINKS [2], if the graph has no missing information.

For the query processing, given q, we can obtain the one-hot encoding of q, denoted by  $\mathbf{x}_q$ . Therefore, given the output node embedding  $\mathbf{Z}$ , the sum of graph shortest-path distances between nodes and the query keywords can be computed with  $\mathbf{x}_q \mathbf{Z}^{\top}$ , and the space complexity is O(NM). It is obvious that Conv-OH cannot deal with the missing information, but it provides some hints to propose more advanced methods.

Conv-PCA. Principal component analysis (PCA) is a classic dimensionality reduction technique
 in multivariate statistical analysis [28]. In order to facilitate data storage and query processing, we
 propose another naive PCA-based method to solve the keyword search problem.

Given the one-hot encoding matrix **X** as the input feature matrix, PCA is able to keep *d* principal components of **X** with  $\mathbf{X}_p = \mathbf{X}\mathbf{U}^{\top}$ , where the rows of  $\mathbf{U} \in \mathbb{R}^{d \times M}$  form an orthogonal basis for the *d* features that are decorrelated [29]. It is worth noting that we can obtain the reconstructed feature matrix **X'** with  $\mathbf{X}' = \mathbf{X}_p \mathbf{U}$ . The learning object is to minimize  $\mathcal{L}_{pca} = ||\mathbf{X}' - \mathbf{X}||_2^2$ , where  $|| \cdot ||_2$ denotes the  $L^2$  norm.

Taking  $\mathbf{X}_p$  as the initial node embedding ( $\mathbf{H}^0 = \mathbf{X}_p$ ), we propose **Conv-PCA** to leverage similar graph convolutional layers of Conv-OH as below:

$$\mathbf{h}_{v}^{l+1} = max\Big(\{\alpha \mathbf{h}_{u}^{l}, \forall u \in \mathcal{N}(v)\} \cup \{\mathbf{h}_{v}^{l}\}\Big),\tag{3}$$

where  $\alpha \in (0, 1)$  is a decay parameter used to estimate the shortest path distances in Eq. (3), since 165 the dimension is reduced from M to d and thus it is difficult to discriminate M keywords within 166 the d dimensions ( $d \ll M$ ). Specifically, a larger cumulative decay corresponds to a larger shortest 167 path distance. This mechanism performs better than directly using PCA in experiments as shown 168 in Section 5.3. For the query processing, given q, we can obtain the query embedding  $h_q = \mathbf{x}_q \mathbf{U}^{\top}$ . 169 Therefore, given the output node embedding Z, the similarity scores between nodes and the query 170 keywords can be computed by  $\mathbf{h}_{q}\mathbf{Z}^{\top}$  with linear space complexity O(dN), where the dimension d is 171 a small constant. 172

Compared with Conv-PCA, Conv-OH utilizes each element of  $h_v$  to record the shortest path distance between keywords and v and cannot reduce the dimension of node embedding. Conv-PCA can more efficiently process the keyword search query and requires less space than Conv-OH, but both of them cannot well handle the missing information in incomplete graphs.

#### 177 4.2 KS-GNN

Based on the prior discussions on Conv-OH and Conv-PCA, we present the desiderata that guide the development of our method for tackling keyword search as follows:

180 **Dimensionality Reduction.** Taking the one-hot encoding matrix X as input, it is difficult to afford 181 the cost of generating an output with size  $N \times M$ . Therefore, the model should be able to reduce the 182 dimensions of the output node embedding M to a lower level.



Figure 3: An illustration of the message passing and aggregation of our KS-GNN model.

183 Key Information Preservation. Some keywords and edges information may be lost in the process of 184 dimensionality reduction, which affects the performance of keyword search. The model should retain 185 as much key information as possible to guarantee results quality.

Adaptive Encoding. When generating node embedding, the model should consider the structure 186 information of the target node centered subgraph and the distribution of keywords on the subgraph, 187 rather than only considering the keywords of the target node. Specifically, to recover the missing 188 keywords information in the incomplete graph, the model should be able to capture latent relationships 189 among different keywords. For instance, a pair of keywords "AI" and "ML" often co-occur on nodes 190 near to each other (e.g., one-hop neighbors). Given a node containing "AI", it is natural to assume that 191 the neighbor of this node is more likely to contain "ML" than the nodes whose one-hop neighbors do 192 not contain "AI". 193

Keyword Frequency Awareness. Based on the scoring function in Section 3, the returned top-k nodes 194 tend to be decided by the query keywords with low frequency compared to the high-frequency ones. 195 Thus, for a given keyword, the number of nodes containing it (we denote this by the keywords' 196 node frequency) can reflect its importance to the query processing, similar to the inverse document 197 frequency (IDF) used in information retrieval. The keyword set of the whole graph can be regarded 198 as a corpus and the keyword set of each node can be regarded as a document. Therefore, the model 199 should encode keywords taking in consideration of their frequencies, which are measured based on 200 the keyword node frequency in the whole graph<sup>2</sup>. 201

Based on these desiderata, we propose an auto-encoder based Keyword Search Graph Neural Network (KS-GNN) for tackling the problem in incomplete graphs. An illustration of the message passing and aggregation mechanism for generating the node embedding  $h_3$  with KS-GNN is provided in Fig. 3.

**Encoder and Decoder.** KS-GNN employs an encoder f to generate low-dimensional node em-205 bedding for dimensionality reduction. Recall that in Conv-PCA, the dimension reduction caused 206 information loss and it is hard to discriminate keywords in the low-dimensional space. To address 207 this issue, for the sake of key information preservation, KS-GNN employs another decoder g which 208 aims to reconstruct the input from embedding space. By training f simultaneously with g, the output 209 embedding of f is able to preserve key information of the input graph. Given the one-hot encoding 210 matrix X as the input, we define  $\mathbf{H} = f(\mathbf{X})$  where  $\mathbf{H} \in \mathbb{R}^{N \times d}$ . For the decoder, it is defined as 211  $\mathbf{X}' = g(\mathbf{H})$  where  $\mathbf{X}' \in \mathbb{R}^{N \times M}$ . 212

In this work, we utilise the multi-layer perceptron (MLP) with a nonlinear activation layer [30] as both non-linear encoder and decoder. It is worth noting that MLP can be replaced by a more complex neural network. Our goal here is to use simple encoder and decoder to show the advantages of the proposed mechanism. In addition, as a conventional learning objective of the auto-encoder, f and gare trained to minimize:

$$\mathcal{L}_1 = \frac{1}{N} ||\mathbf{X}' - \mathbf{X}||_2^2.$$
(4)

It is worth noting that the representation PCA learns is essentially the same as that learned by a basic linear auto-encoder, but the encoder f here is not required to generate embedding based on the primary components.

<sup>220</sup> printary components.

<sup>&</sup>lt;sup>2</sup>The multi-occurrence of a keyword on one node only contribute 1 to this keyword's node frequency.

**Message Passing and Aggregation.** The message passing and aggregation mechanisms for both 221 Conv-OH and Conv-PCA are based on the orthogonal basis and decorrelated features. However, 222 in KS-GNN, its encoder f transforms the input without any basis, and thus we cannot simply 223 apply the  $\max(\cdot)$  function to capture the information of the nearest keywords on a node. Thanks 224 to the reconstruction ability of the decoder q, we can utilise q to reconstruct the M-dimension 225 encoding during the message passing and then using  $\max(\cdot)$  at this step. Moreover, if the learned 226 node embedding contains the latent information of the missing keyword of the incomplete graph, q 227 can also help recover the missing keywords during the reconstruction. Formally, given the output 228 node embedding  $\mathbf{H}^l$  of *l*-th layer, we have: 229

$$\mathbf{h}_{v}^{l+1} = f\bigg(max\Big(\{\alpha g(\mathbf{h}_{u}^{l}), \forall u \in \mathcal{N}(v)\} \cup \{g(\mathbf{h}_{v}^{l})\}\Big)\bigg),\tag{5}$$

where  $\alpha \in (0, 1)$  is a decay parameter that is the same as the one in Eq. (3). As Eq. (5) shows, the *M*dimension embedding will only be generated by *g* during the message passing and aggregation, while the hidden node embedding and the final output node embedding are both *d* dimensions. Therefore, it meets the requirement of *dimensionality reduction*. The message passing and aggregation can be processed in parallel, and the time-complexity is acceptable. It is worth noting that in incomplete graphs, due to the mechanism of message passing and aggregation, nodes without any keyword information can still be embedded accordingly.

**Subgraph keywords-based Node Similarity.** To realize *adaptive encoding*, we propose to train 237 KS-GNN by a triplet siamese network [31] with a triplet loss [32] according to the subgraph 238 keywords-based node similarity, which enables KS-GNN to capture the latent missing keyword and 239 edge information on incomplete graphs. For a node v, we consider the subgraph  $SG_v$  containing all 240 the neighbors of v within k hops for measuring node similarity. The one-hot encoding of this subgraph 241  $SG_v$  is denoted by  $\mathbf{x}_{SG_v}$ . For instance in G' shown in Fig. 1, given k = 1, the 1-hop subgraph of  $v_5$ 242 contains keywords  $\{a, c, d, e, f\}$  with the corresponding one-hot encoding  $\mathbf{x}_{SG'_5} = (1, 0, 1, 1, 1, 1)$ . 243 Similarly, the one-hot encoding of the subgraphs around  $v_4$  and  $v_6$  are  $\mathbf{x}_{SG'_4} = (1, 1, 1, 0, 0, 1)$  and 244  $\mathbf{x}_{SG'_6} = (1, 1, 1, 1, 1, 1)$ , respectively. Therefore, taking the dot product of embedding as the simi-245 larity scoring function, we can compare the similarity between (v5, v4) and (v5, v6) by comparing  $\mathbf{x}_{SG'_5}\mathbf{x}_{SG'_4}^{\top} = 3$  and  $\mathbf{x}_{SG'_5}\mathbf{x}_{SG'_6}^{\top} = 5$ . Specifically,  $\mathbf{x}_{SG'_5}\mathbf{x}_{SG'_4}^{\top} < \mathbf{x}_{SG'_5}\mathbf{x}_{SG'_6}^{\top}$  indicates (v5, v6) are more similar than (v5, v4) in G'. 246 247 248

Given G', the KS-GNN model denoted by  $\phi$ , and a sampled batch of triplets  $\mathcal{T} = \{t_1, t_2, ...t_n\} = \{(v_{o1}, v_{p1}, v_{q1}), (v_{o2}, v_{p2}, v_{q2}), ..., (v_{on}, v_{pn}, v_{qn})\}$ , KS-GNN is trained to minimise:

$$\mathcal{L}_{2} = \frac{1}{|\mathcal{T}|} \sum_{t_{i} \in \mathcal{T}} max \bigg( m - \mathbb{1}(t_{i}) \Big( \phi(\mathbf{X})_{o_{i}} \phi(\mathbf{X})_{p_{i}}^{\top} - \phi(\mathbf{X})_{o_{i}} \phi(\mathbf{X})_{q_{i}}^{\top} \Big), 0 \bigg),$$
(6)

where *m* is a margin hyper-parameter of the hinge loss,  $1(\cdot)$  denotes an indicator function that  $1(t_i)$ returns 1 if  $(v_{oi}, v_{pi})$  are more similar than  $(v_{oi}, v_{qi})$  and -1 otherwise. Thus, KS-GNN can learn the structure and keyword information from the subgraphs involved.

For large-scale datasets, it might be time-consuming to compute  $1(t_i)$ . In this case, it is acceptable to intuitively sample  $\mathcal{T}$  based on the links. For example, for a sampled node  $v_{o_i}$ ,  $v_{p_i}$  can be sampled from the 1-hop neighbors of  $v_{o_i}$ , and  $v_{q_i}$  can be negatively sampled from unconnected nodes of  $v_{o_i}$ , thereby setting  $1(t_i)$  to 1. Eq. (6) still takes both the graph structural information and the keyword distribution into account by feeding the one-hot encoding of subgraph keywords to KS-GNN. In addition, minimizing Eq. (6) helps generate similar adaptive embedding for the keywords which co-occur commonly.

**Keyword Frequency-based Regularization.** Intuitively, if a keyword appears on many nodes, it is regarded as less important than the keyword which appears on fewer nodes for query processing. Therefore, in this work, we consider the keyword node frequency, denoted by  $c_i$ , that indicates the number of nodes containing keyword  $w_i$ . For instance, in G' shown in Fig. 1,  $c_1 = 4$  and  $c_2 = 2$  for keywords a and b, respectively. We propose to enhance the model's keyword frequency awareness with a regularization that minimizes:

$$\mathcal{L}_3 = \frac{1}{M} \sum_{w_i \in \mathcal{W}} c_i ||f(\mathbf{I}_i)||_2, \tag{7}$$

where I denotes an  $M \times M$  identity matrix, and  $I_i$  denotes the *i*-th row of I. Feeding  $I_i$  in f can return the representation of keyword  $w_i$ , and minimizing Eq. (7) aims to differentiate the lengths of keyword embeddings according to their keyword frequencies, thereby being aware of keyword

270 frequency.

271 To train KS-GNN, the final learning objective is to minimize:

$$\mathcal{L} = \lambda_1 \mathcal{L}_1 + \lambda_2 \mathcal{L}_2 + \lambda_3 \mathcal{L}_3, \tag{8}$$

where  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are hyper-parameters. By minimizing Eq. (8), we can optimize KS-GNN to generate informative node embedding which can capture the latent representation of missing keywords and edges. The superiority of the proposed KS-GNN is validated in Section 5.3.

**Query Processing.** To process query q, given the one-hot encoding of q as  $\mathbf{x}_q$ , the trained encoder f and the learned node embedding  $\mathbf{Z}$ , we can compute the similarity between the nodes and query with  $\mathbf{s}_q = f(\mathbf{x}_q)\mathbf{Z}^{\top}$ , and the top-k answers can be found with the largest scores in  $\mathbf{s}_q$ . In addition, the space complexity of computing query processing is O(dN).

## 279 **5 Experiments**

In this section, we evaluate the performance of our proposed approach, KS-GNN, on four real-world datasets, including citation networks (**CiteSeer**), co-purchase networks (**Video & Toy**) and co-author networks (**DBLP**). The details of datasets, additional experimental results and analysis can be found in the supplementary materials.

#### 284 **5.1 Baseline Methods**

We compare our model against five baseline methods, including a state-of-the-art deep learning based missing-data completion GNN model. More details on the baseline models are provided in the supplementary materials.

- **GraphSAGE** [7] is a representative GNN-based graph embedding method. We add an MLP encoder for GraphSAGE to address the keyword search problem for GraphSAGE.
- BLINK+SAT firstly predicts and completes the missing keywords and edges with a stateof-the-art missing-data completion GNN model SAT [6] and then utilises BLINK [2] to process keyword search on the new graph.
- **PCA** is based on the classic dimensionality reduction technique [28].
- **Conv-PCA** is a naive method proposed in Section 4.1.
- Conv-rPCA is a variant of Conv-PCA that leverages U to reconstruct M-dimension embedding from  $\mathbf{h}_v$ .

#### 297 5.2 Experimental Setup

In our experiments, we compare the proposed method with baseline methods for keyword search 298 tasks in two kinds of graphs: (1) the graphs with only missing keywords; (2) the graphs with both 299 missing keywords and edges. For each dataset, to simulate a real-world scenario and quantitatively 300 control the ratios of missing information, we process the original datasets with two steps: (1) hide the 301 keywords of randomly sampled nodes with a predefined proportion (denoted by  $r_w$ ) in the graph; 302 (2) randomly hide a proportion (denoted by  $r_e$ ) of the edges in the graph. Let  $n_q = |q|$  denote the 303 number of words in the query q, we randomly sample 100 queries as the test set for each value of  $n_q$ 304 ranging from 3 to 9 with a step of 2. 305

In addition, in each incomplete graph, the validation set consists of 100 randomly generated queries 306 with ground truth answers. We tune the hyper-parameters of compared methods with the grid search 307 algorithm on the validation set, more details can be found in the Appendix. In terms of the evaluation 308 metric, we use Hits@K, which is a common ranking metric that counts the ratio of positive edges that 309 are ranked at the K-th place or above. The ground truth is the top-K answers retrieved by BLINK on 310 the original graph for each dataset. Specifically, we report Hits@100, and more experimental results 311 (Hits@10 and Hits@50) can be found in the appendix. Moreover, for each experiment, we conduct 312 10 runs and report the average Hits scores. 313

Datasets	$r_w$	0.3					0	.5		0.7				
	$n_q$	3	5	7	9	3	5	7	9	3	5	7	9	
Citeseer	GraphSAGE	1.72	5.38	6.72	3.91	9.33	4.15	2.77	4.99	9.68	7.26	8.29	6.75	
	BLINK+SAT	9.8	10.88	12.87	14.49	9.86	8.36	6.12	9.67	3.57	1.37	1.51	2.07	
	PCA	13.12	8.41	7.04	6.05	9.86	7.34	6.80	5.66	8.40	6.83	6.07	5.46	
	Conv-PCA	8.10	7.20	7.52	7.87	10.93	8.15	9.38	5.23	8.67	6.45	6.93	4.83	
	Conv-rPCA	24.66	28.91	33.79	34.92	24.99	26.72	29.85	32.52	22.37	27.15	31.11	34.42	
	KS-GNN	31.21	42.13	39.89	41.53	33.19	41.73	41.79	41.24	32.56	42.74	43.71	42.45	
Video	GraphSAGE	0.49	0.30	0.06	0.03	0.44	0.14	0.00	0.05	0.34	0.35	0.21	0.16	
	BLINK+SAT	10.21	9.86	10.99	14.87	8.55	6.92	8.63	5.82	1.18	1.15	4.38	3.35	
	PCA	1.54	0.91	0.55	0.61	1.71	0.72	0.71	0.57	1.66	0.95	0.66	0.55	
	Conv-PCA	1.81	2.46	1.58	2.38	2.43	1.49	1.66	2.54	2.42	2.37	2.80	3.37	
	Conv-rPCA	10.19	12.23	16.15	21.37	11.26	15.62	19.51	23.87	10.66	15.57	19.17	25.36	
	KS-GNN	21.43	26.63	22.92	39.51	22.54	22.57	30.41	38.63	21.01	16.48	22.01	39.19	
Тоу	GraphSAGE	0.09	2.15	2.09	10.07	0.03	1.35	0.50	13.13	0.06	0.30	0.00	3.22	
	BLINK+SAT	11.02	9.73	7.71	11.28	8.97	7.91	7.73	6.14	0.63	1.23	1.15	1.77	
	PCA	1.85	0.47	0.66	0.43	1.40	0.42	0.43	0.39	1.23	0.42	0.45	0.34	
	Conv-PCA	15.37	15.92	17.96	14.60	11.92	11.37	16.35	12.95	12.48	11.64	13.44	11.81	
	Conv-rPCA	23.78	24.14	22.95	22.09	24.57	23.82	26.38	22.30	15.02	20.03	21.24	21.67	
	KS-GNN	25.42	26.41	28.84	28.51	29.26	34.03	30.12	25.72	18.82	25.29	35.15	23.16	
DBLP	GraphSAGE	0.05	0.27	0.06	0.23	0.00	0.06	1.06	1.14	1.05	0.28	0.18	1.40	
	BLINK+SAT	8.37	9.97	8.05	9.89	3.91	4.01	3.23	4.56	4.75	1.93	3.53	3.29	
	PCA	3.22	2.95	2.22	2.08	3.22	2.65	2.27	1.90	2.83	2.30	1.76	1.74	
	Conv-PCA	5.51	5.96	5.19	6.28	4.32	7.36	5.51	7.67	3.18	4.37	1.21	2.58	
	Conv-rPCA	14.40	24.30	25.82	24.96	12.59	18.51	22.93	22.62	9.95	19.97	20.87	23.39	
	KS-GNN	24.41	34.67	39.98	43.70	22.21	31.56	39.97	40.99	20.44	28.84	33.96	36.73	

Table 1: Method performance by Hits@100 (%) in graphs with only missing keywords.

#### 314 5.3 Performance of Keyword Search

Table 1 shows the comparison results in graphs with  $r_w$  adjusted from 0.3 to 0.7 and  $n_q$  adjusted 315 from 3 to 9. As shown in the table, KS-GNN significantly outperforms the baselines, and changing 316  $r_w$  will not affect its performance. Moreover, the performance of KS-GNN is better when more 317 query keywords are given. As for the baselines, BLINK+SAT cannot maintain good performance 318 when many keywords are missing. Compared with Conv-PCA, Conv-rPCA can address the keyword 319 search in incomplete graphs much more effectively due to the proposed novel message passing 320 and aggregation mechanism. Although Conv-rPCA has the same message-related mechanism as 321 KS-GNN, KS-GNN always performs better. This reveals that the proposed learning objective and 322 auto-encoder-based model are able to enhance the ability of representation learning. Since PCA 323 focuses on each single node, it performs well when the query keywords are located on the same 324 node. However, when  $n_q$  increases, the query keywords tend to be located on different nodes, and 325 326 the performance of PCA therefore decreases because it cannot gather neighbor information. By contrast, although GraphSAGE can aggregate information from neighbors, it sometimes performs 327 worse than PCA, because only utilizing the max-pooling operator during the message aggregation 328 cannot well distinguish the information from each unique keyword. This can be proved by that 329 Conv-PCA performs better than both PCA and GraphSAGE, which also indicates the superiority of 330 our proposed encoder and decoder-based message passing and aggregation mechanism. 331

Table 2 presents the results of the comparison in graphs with both missing keywords and edges, where  $r_e$  is set to 0.3 and  $r_w$  is adjusted from 0.3 to 0.7. As the table shows, KS-GNN significantly outperforms all the compared baseline methods, since only KS-GNN can learn the adaptive embedding and structural information from the incomplete graph with missing keywords and edges. Compared with Table 1, Table 2 shows that edge missing has no significant effect on the performance of KS-GNN, while the performance of the baseline methods decreases significantly with the increase of the missing edges. This proves the robustness of our proposed KS-GNN.

#### 339 5.4 Analysis of Keyword Frequency Awareness

Datasets	$r_w$		0	.3			0	.5		0.7			
	$n_q$	3	5	7	9	3	5	7	9	3	5	7	9
Citeseer	GraphSAGE	6.95	4.75	7.29	4.03	10.70	3.97	2.00	3.34	6.62	7.88	9.54	3.66
	BLINK+SAT	6.67	7.33	8.27	8.47	4.24	5.23	5.85	4.92	1.04	2.64	2.91	2.21
	PCA	11.24	8.41	7.04	6.05	9.86	7.34	6.80	5.66	8.40	6.83	6.07	5.46
	Conv-PCA	10.27	8.17	7.40	8.31	11.83	8.47	10.15	5.33	9.12	7.39	8.32	4.58
	Conv-rPCA	25.16	28.95	31.00	33.22	18.89	25.93	28.25	30.52	25.24	25.33	30.42	30.06
	KS-GNN	29.97	38.12	40.32	41.63	30.17	38.56	36.61	41.94	30.61	40.29	40.65	40.10
Video	GraphSAGE	0.09	0.21	0.02	0.00	0.26	0.05	0.00	0.04	1.65	1.65	2.22	1.29
	BLINK+SAT	1.67	1.85	2.48	1.44	0.08	0.99	4.96	2.97	2.19	1.77	0.78	1.21
	PCA	1.54	0.91	0.55	0.61	1.71	0.72	0.71	0.57	1.66	0.95	0.66	0.55
	Conv-PCA	1.05	1.59	0.83	2.01	1.43	0.81	0.74	1.23	1.25	1.25	1.31	1.38
	Conv-rPCA	3.82	4.13	4.88	7.13	3.96	4.52	5.33	6.11	3.64	4.58	5.17	6.55
	KS-GNN	12.81	8.34	12.88	21.86	10.01	7.68	7.18	20.30	10.34	10.31	13.90	19.03
Тоу	GraphSAGE	0.00	0.14	0.74	1.24	0.01	0.24	0.46	4.77	0.00	0.21	0.23	0.71
	BLINK+SAT	6.40	4.32	2.59	6.79	3.53	2.68	5.92	2.16	1.87	1.65	0.94	0.54
	PCA	1.85	0.47	0.66	0.43	1.40	0.42	0.43	0.39	1.23	0.42	0.45	0.34
	Conv-PCA	7.92	7.00	7.40	7.04	4.66	3.37	5.06	4.13	4.76	4.49	5.37	4.55
	Conv-rPCA	10.12	9.19	11.36	12.03	7.51	11.04	8.54	11.35	7.83	7.95	5.49	8.14
	KS-GNN	12.83	12.28	12.66	13.66	12.09	11.96	10.88	15.80	8.27	12.97	8.92	9.34
DBLP	GraphSAGE	0.21	0.27	0.19	0.47	1.98	3.14	1.11	3.57	1.42	0.92	1.04	1.47
	BLINK+SAT	9.01	9.91	6.79	8.54	5.56	4.86	4.35	2.71	1.7	0.36	0.68	0.03
	РСА	3.22	2.95	2.22	2.08	3.22	2.65	2.27	1.90	2.83	2.30	1.76	1.74
	Conv-PCA	4.49	5.48	4.89	5.81	3.35	5.92	4.46	6.17	3.23	4.26	1.32	2.36
	Conv-rPCA	5.11	12.55	13.95	14.39	11.30	16.76	16.06	4.66	4.06	10.10	6.41	16.68
	KS-GNN	22.07	26.14	33.67	36.09	17.26	25.20	28.26	32.14	17.12	24.69	27.18	28.14

Table 2: Method performance by Hits@100 (%) in graphs with both missing keywords and edges ( $r_e = 0.3$ ).

As discussed above in Section 4.2, we propose 340 a novel learning objective for training KS-GNN 341 that aims to enhance its ability of keyword fre-342 quency awareness. Therefore, in the incomplete 343 graph with  $r_w = 0.3$  and  $r_e = 0$ , we conduct ex-344 periments which show the relation between the 345 keyword frequency  $c_i$  and the length of keyword 346 embedding  $||f(\mathbf{I}_i)||_2$ . We compare the results by 347 setting  $\lambda_3$  to 1 or 0, which indicates whether to 348 minimize Eq. (7) or not. 349



Figure 4: Comparison of the ability of keyword frequency awareness by whether using  $\mathcal{L}_3$  or not.

As the figure shows, by minimizing Eq.(7), KS-GNN can significantly learn the keyword frequency awareness, which is reflected by the length of keyword embedding. It is presented that the keywords with high frequencies turn to be less important than before utilizing Eq. (7). Because compared with the long keyword embedding, shorter keyword embedding tends to be ignored during the query process. It is also interesting to notice that the lengths of some low-frequency keywords decrease. This is exactly what we expect since there are many low-frequency keywords in the graph, therefore it is meaningful to distinguish them according to their importance.

## 357 6 Conclusion

Keyword search in graphs is an important problem with many applications such as network analysis 358 and recommendation. The keywords and edges in graphs might be lost or incomplete due to some 359 reasons in real-world applications, such as storage limitation or privacy issues. In this paper, we study 360 the keyword search problem in incomplete graphs and propose a novel auto-encoder and GNN-based 361 method, KS-GNN. Compared to existing methods, KS-GNN is able to address the problem when 362 some nodes have missing keywords or some edges are missing in the input graphs. The results of 363 extensive experiments on real-world datasets reveal that KS-GNN significantly outperforms the 364 state-of-the-art baseline methods on the incomplete graph keyword search task. 365

## 366 **References**

- [1] Haixun Wang and Charu C. Aggarwal. A survey of algorithms for keyword search on graph
   data. In *Managing and Mining Graph Data*, volume 40 of *Advances in Database Systems*, pages
   249–273. Springer, 2010.
- [2] Hao He, Haixun Wang, Jun Yang, and Philip S Yu. Blinks: ranked keyword searches on graphs.
   In *SIGMOD*, pages 305–316, 2007.
- [3] Wangchao Le, Feifei Li, Anastasios Kementsietsidis, and Songyun Duan. Scalable keyword
   search on large rdf data. *IEEE Transactions on knowledge and data engineering*, 26(11):2774–2788, 2014.
- [4] Jieming Shi, Dingming Wu, and Nikos Mamoulis. Top-k relevant semantic place retrieval on
   spatial rdf data. In *SIGMOD*, pages 1977–1990, 2016.
- [5] Dejian Yang, Senzhang Wang, Chaozhuo Li, Xiaoming Zhang, and Zhoujun Li. From properties to links: Deep network embedding on incomplete graphs. In *CIKM*, pages 367–376, 2017.
- [6] Xu Chen, Siheng Chen, Jiangchao Yao, Huangjie Zheng, Ya Zhang, and Ivor W Tsang. Learning
   on attribute-missing graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
   2020.
- [7] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Proc. of NeurIPS*, 2017.
- [8] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn:
   An efficient algorithm for training deep and large graph convolutional networks. In *SIGKDD*,
   pages 257–266, 2019.
- [9] S. Agrawal, S. Chaudhuri, and G. Das. Dbxplorer: a system for keyword-based search over
   relational databases. In *ICDE*, pages 5–16, 2002.
- [10] Gaurav Bhalotia, Arvind Hulgeri, Charuta Nakhe, Soumen Chakrabarti, and Shashank Sudar shan. Keyword searching and browsing in databases using banks. In *ICDE*, pages 431–440.
   IEEE, 2002.
- [11] Mehdi Kargar and Aijun An. Keyword search in graphs: Finding r-cliques. *Proceedings of the VLDB Endowment*, 4(10):681–692, 2011.
- [12] Ziyang Liu, Chong Wang, and Yi Chen. Keyword search on temporal graphs. *IEEE Trans. Knowl. Data Eng.*, 29(8):1667–1680, 2017.
- [13] Ye Yuan, Guoren Wang, Lei Chen, and Haixun Wang. Efficient keyword search on uncertain
   graph data. *IEEE Trans. Knowl. Data Eng.*, 25(12):2767–2779, 2013.
- [14] Yueji Yang, Divyakant Agrawal, H. V. Jagadish, Anthony K. H. Tung, and Shuang Wu. An
   efficient parallel keyword search engine on knowledge graphs. In *ICDE*, pages 338–349. IEEE,
   2019.
- [15] Shady Elbassuoni and Roi Blanco. Keyword search over RDF graphs. In *CIKM*, pages 237–242.
   ACM, 2011.
- [16] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally
   connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.
- [17] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks
   on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29:3844–3852, 2016.
- [18] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
   networks. In *Proc. of ICLR*, 2017.
- [19] David I Shuman, Sunil K Narang, Pascal Frossard, Antonio Ortega, and Pierre Vandergheynst.
   The emerging field of signal processing on graphs: Extending high-dimensional data analysis to
   networks and other irregular domains. *IEEE signal processing magazine*, 30(3):83–98, 2013.

- [20] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
   networks? In *ICLR*, 2018.
- [21] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
   Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- 417 [22] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *NeurIPS Workshop*, 2016.
- [23] Jiaxuan You, Xiaobai Ma, Daisy Yi Ding, Mykel Kochenderfer, and Jure Leskovec. Handling
   missing data with graph representation learning. *arXiv preprint arXiv:2010.16418*, 2020.
- <sup>420</sup> [24] Indro Spinelli, Simone Scardapane, and Aurelio Uncini. Missing data imputation with <sup>421</sup> adversarially-trained graph convolutional networks. *Neural Networks*, 129:249–260, 2020.
- [25] Bhavana Bharat Dalvi, Meghana Kshirsagar, and S Sudarshan. Keyword search on external
   memory data graphs. *Proceedings of the VLDB Endowment*, 1(1):1189–1204, 2008.
- [26] Gang Gou and Rada Chirkova. Efficient algorithms for exact ranked twig-pattern matching over
   graphs. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 581–594, 2008.
- [27] Thanh Tran, Haofen Wang, Sebastian Rudolph, and Philipp Cimiano. Top-k exploration of
   query candidates for efficient keyword search on graph-shaped (rdf) data. In 2009 IEEE 25th
   International Conference on Data Engineering, pages 405–416. IEEE, 2009.
- [28] Jan J Gerbrands. On the relationships between svd, klt and pca. *Pattern recognition*, 14(1-6):375–381, 1981.
- [29] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and
   new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–
   1828, 2013.
- [30] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- [31] Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *International* workshop on similarity-based pattern recognition, pages 84–92. Springer, 2015.
- [32] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for
   face recognition and clustering. In *Proc. of CVPR*, pages 815–823, 2015.

## 440 Checklist

446

447

448

449

451

452

453

454

455

- 441 1. For all authors...
- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes] This work focuses on handling keyword search problem with missing information.
  (b) Did you describe the limitations of your work? [Yes] See section 6. This research needs
  - (b) Did you describe the limitations of your work? [Yes] See section 6. This research needs to be extended for larger-scale datasets.
  - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
    - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 450 2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
  - (b) Did you include complete proofs of all theoretical results? [N/A]
  - 3. If you ran experiments...
    - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See sections starting from section 5.2.

458 459	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] See section 5.2.
460 461	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix C.1.
462	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
463 464	<ul><li>(a) If your work uses existing assets, did you cite the creators? [Yes] See Section ??.</li><li>(b) Did you mention the license of the assets? [Yes]</li></ul>
465	(c) Did you include any new assets either in the supplemental material or as a URL? [No]
466 467	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
468 469	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
470	5. If you used crowdsourcing or conducted research with human subjects
471 472	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
473 474	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
475 476	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]