

# InducT-GCN: Inductive Graph Convolutional Networks for Text Classification

Kunze Wang, Soyeon Caren Han\*, Josiah Poon  
School of Computer Science, The University of Sydney  
kwan4418@uni.sydney.edu.au, {caren.han, josiah.poon}@sydney.edu.au

**Abstract**—Text classification aims to assign labels to textual units by making use of global information. Recent studies have applied graph neural network (GNN) to capture the global word co-occurrence in a corpus. Existing approaches require that all the nodes (training and test) in a graph are present during training, which are transductive and do not naturally generalise to unseen nodes. To make those models *inductive*, they use extra resources, like pretrained word embedding. However, high-quality resource is not always available and hard to train. Under the extreme settings with no extra resource and limited amount of training set, can we still learn an inductive graph-based text classification model? In this paper, we introduce a novel inductive graph-based text classification framework, **InducT-GCN** (InducTive Graph Convolutional Networks for Text classification). Compared to transductive models that require test documents in training, we construct a graph based on the statistics of training documents only and represent document vectors with a weighted sum of word vectors. We then conduct one-directional GCN propagation during testing. Across five text classification benchmarks, our InducT-GCN outperformed state-of-the-art methods that are either transductive in nature or pre-trained additional resources. We also conducted scalability testing by gradually increasing the data size and revealed that our InducT-GCN can reduce the time and space complexity. The code is available on: <https://github.com/usydnlp/InducTGCN>.

## I. INTRODUCTION

Text classification is one of the most fundamental natural language processing research tasks, including topic classification, news categorisation, and sentiment analysis. The aim of the text classification is to classify/categorise textual documents into the predefined classes. The high-dimensional textual input is assigned to the output class with binary or multi-class classification. Note that we only consider a single-label text classification problem in this research.

Many recent text classification studies have focused on learning text representations using sequence-based learning models, such as convolutional neural networks (CNN) [1] or recurrent neural networks (RNN) /long short term memory (LSTM) [2]. The CNN/RNN-based models focus on the locality and sequence of text and mainly aim to detect semantic and syntactic information in local consecutive word sequences. It tends to neglect global word co-occurrence in a corpus and ignore non-consecutive and long-distance semantic information [3]. However, those models need a relatively large training set to perform better. Still, most real-world cases (e.g., specific domain or some low resource languages) have

a very limited amount of training set (limited labeled data). Recently, pre-trained models, like BERT [4], RoBERTa [5], have achieved state-of-the-art performance on several NLP tasks with limited amount of training data. However, those models require much computation and external resources for pre-training, which are not always available.

[6] proposed TextGCN and performed well, especially when the percentage of training data is low without using any external resources and with low computation costs. It is an initial text GNN framework, which conducts a straightforward manner of graph construction and applies a GCN-learning [7] to deal with complex structured textual data and prioritise global feature exploitation. More recent studies [8]–[10] utilise more contextual information or optimising the computation.

However, most graph models are intrinsically transductive. The learned node representations/embeddings for words/documents are not naturally generalisable to unseen words/documents, making it challenging to apply in the real world. The transductive nature of these graph-based learning models requires relatively large computational space when the corpus size is large. Therefore, an inductive model is needed. To expand a transductive graph-based text classification into an inductive model, we mainly consider the following three requirements: 1)The inductive learning model must not include any test set information during the training. 2)The inductive model must not re-train the model on the whole new graph when it learns a new sample. 3)We use corpus-level graph-based text classification to make an inductive model. It nicely covers the benefit of GNN, which captures the complex global structure of the whole corpus and prioritises global feature exploitation. In this paper, we propose a novel inductive graph-based text classification framework, called **InducT-GCN** (InducTive Graph Convolutional Networks for Text classification). We introduce a new inductive graph framework of graph construction, learning, and testing, and it can expand to any transductive GCN-based text classification model. The paper includes the following contributions:

- To the best of our knowledge, we introduce the first inductive corpus-level GCN-based text classification framework without using extra resources.
- We compare our InducT-GCN on five benchmark datasets under the limited labeled data settings. InducT-GCN outperforms on four of them, even beating some transductive baselines integrated using external resources.
- We introduce a new way to make transductive GCN-based

\* Corresponding author (Caren.Han@sydney.edu.au)

text classification models *inductive*, which improves the performance and reduces the time and space complexity.

## II. RELATED WORK

### A. Graph Neural Networks

Graph Neural Network (GNN)s [7] have been effective at tasks to have rich relational structure and can preserve global structure information of a graph in graph embeddings by aggregating first-order neighbourhood information. [7] introduced Graph Convolutional Networks (GCN) on transductive classification tasks. GraphSage [11], and FastGCN [12] tailored GCNs on inductive representation learning framework with sampling methods. Graph Attention Networks (GAT) [13] applied the Attention to specify different weights to different nodes in a neighbourhood. More recent GCN studies for transductive and inductive frameworks have been proposed. For transductive-based GCN, SGC [8] was introduced to reduce the complexity and S<sup>2</sup>GC [10] was proposed to solve over-smoothing problems. Some inductive-based models, DeepGL [14] and TGAT [15], were introduced to cover different graph tasks, including transfer learning and topology learning.

### B. Text Classification Using GNN

GNNs have received attention in various NLP tasks [6], [16]–[20], including text classification. The GNN-based text classification models can be categorised into two types, Document-level and Corpus-level approaches.

**Document-level GNN in Text Classification** Several graph-based text classification models build a graph for each document using words as nodes [3], [21]–[25]. Word nodes are represented by external resources, pre-trained embedding, such as Word2vec [26], and Glove [27]. The edges are built either using word co-occurrence information [3], [25] or simply connecting consecutive words in a sentence [24]. Hence, they do not consider explicit global structure information of a corpus/entire dataset during their model training/learning.

**Corpus-level GNN in Text Classification** TextGCN [6] was proposed to build a graph for the entire text corpus with documents and words as nodes. Hence, it captures global information of an entire corpus and conduct node(document) classification. SGC [8] and S<sup>2</sup>GC [10] constructed a graph as TextGCN, but proposed different information propagation approaches. Both TensorGCN [9] and TextGTL [28] proposed three graphs to cover three different aspects. Note that all three graphs are based on an entire corpus and use the same propagation as GCN [7]. TG-Transformer [29] applied transformer with pretrained GloVe embeddings to the TextGCN, and BERTGCN [30] applied BERT embedding to the TextGCN. All the above models are transductive-based approaches as GCN [7]. However, our model Induct-GCN, an inductive graph-based text classification framework, constructs a corpus-level graph but adopts the nature of inductive learning to generalise to unseen nodes naturally.

## III. INDUCT-GCN

We propose an Inductive Graph Convolutional Network (GCN) for text classification, named ‘Induct-GCN’, which can be an extension of the traditional transductive GCN-based text classification models. We adopt the traditional transductive GCN-based text classification models, including TextGCN [6] and SGC [8], and focus on expanding those models to efficient inductive learning models. This section demonstrates the proposed inductive learning components applied to TextGCN.

### A. Revisit TextGCN

TextGCN is a GCN-based text classification model that uses a large text graph based on the whole corpus. To understand the concept properly, we first explore the GCN process.

**Graph Convolutional Networks(GCN)** Formally, considering a graph  $G = (V, E, A)$ , where  $V(|V| = n)$  is a set of nodes,  $E$  is a set of edges, and  $A \in R^{n \times n}$  is an adjacency matrix representing the edge values between nodes. The propagation rule of each hidden layer is:

$$H^{(l+1)} = f(H^{(l)}, A) = \sigma(\tilde{A}H^{(l)}W^{(l)}) \quad (1)$$

where  $\tilde{A} = D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$  is a normalized symmetric matrix for  $A$  and  $D_{ii} = \sum_j A_{ij}$  as a degree matrix of adjacency matrix  $A$ .  $H^{(l)}$  is  $l_{th}$  hidden layer input and  $W^{(l)}$  is the weight to be learned in this layer.  $\sigma$  is an activation function, e.g. ReLU:  $\sigma(x) = \max(0, x)$ .

**TextGCN** Followed by the GCN [7], TextGCN constructs a large corpus-level graph but with textual information, documents and words as nodes so it can model the global word-document co-occurrence. The constructed graph includes documents and words nodes from training sets and test sets. TextGCN aims to model the global word-document occurrence with two major edges: 1) word-word edge: calculated by co-occurrence information point-wise mutual information(PMI) [31], 2) document-word edge: TF-IDF. One-hot vectors are fed into a two-layer GCN model to jointly learn the embedding for the documents and words. The representations on the document nodes in the training set train the classification model while those in the test set are used for prediction.

### B. Transductive and Inductive Nature

This section discusses the nature of transductive and inductive GCN learning for text classification and what inductive learning aspect we would like to explore. Most GCN text classification models, including TextGCN [6], SGC [8], or S<sup>2</sup>GC [10], are inherently transductive by using the whole corpus, including training set and test set all time.

To expand those transductive models into an inductive learning nature, we fundamentally improve two aspects as follows. First, the transductive GCN-based text classification models include documents from the training set and the test set when constructing a whole-corpus-based textual graph for GCN learning. Hence, the learned GCN model will be influenced/generalised by word/document information in the test set, which is supposed to be unseen nodes. Our inductive GCN-based text classification model constructs a graph with

only training document information but does not consider any information from the test sets. We focus on generalising to unseen nodes and aligning newly observed subgraphs to the node that the model has already optimised on.

Secondly, the transductive models learn the embedding for  $V_{train}$ ,  $V_{test}$ ,  $V_{word}$  simultaneously by using one-hot input vectors  $H^{(0)} \in R^{n \times n}$ . For any new test sample, the embedding should be re-learned by re-training the model on the new graph. In this case, the re-learning/re-training process does not perfectly fulfil the effective generalisation to unseen nodes. Therefore, we develop a new graph construction and training/testing solution for inductive learning instead of re-learning or re-training.

### C. InducT-GCN Graph Construction

1) *Graph Nodes*: Our inductive GCN-based text classification model, InducT-GCN, strictly do not consider any information or statistics from the test set, which is supposed to be unseen nodes. Instead, we construct the nodes only with training document information. Consider a set of nodes  $V = \{V_{train}, V_{word}\}$  and the  $V_{word}$  are the unique words in the training documents. To define input vector  $H^{(0)}$  for graph nodes in the InducT-GCN graph, we consider two requirements: (1) During the propagation phase, the graph is considered as a homogeneous graph, which means all the nodes are considered as the same type without checking whether they are word nodes or document nodes. Then all the input vectors for document nodes and word nodes should align with each other. (2) Our InducT-GCN must not use one-hot vectors for representing document nodes to avoid learning any representation on testing documents during training.

With this in mind, we propose a new document representation by focusing on the nature of our proposed inductive learning idea. InducT-GCN generates document node representations based on its word node vectors for the proper alignment between word and document representation. We use a weighted average of word vectors to construct document nodes vectors, and the key idea of this construction is applying TF-IDF weights. Formally, one-hot vectors are used for representing word nodes vectors  $H_i^{(0)}, \forall i \in V_{word}$ . For representing training documents node vectors  $H_i^{(0)}, \forall i \in V_{train}$ , we use TF-IDF vectors. The values for each dimension is TF-IDF values for the corresponding word in that specific document:  $H_{ij}^{(0)} = \text{TF-IDF}(i, j)$  where  $i$  and  $j$  are document and word, respectively. Figure 1 shows an example of  $H^{(0)}$ .

2) *Graph Edges*: We focus on expanding corpus-level transductive GCN-based text classification to inductive learning and select TextGCN [6] as one of its kind. Like TextGCN, we define two-edge types for the InducT-GCN graph: 1) Word-Word with PMI and 2) Word-Doc edges with TF-IDF. Note that each node also connects to itself. PMI is calculated based on the co-occurrence of a pair of words in a slicing window. Formally, it is calculated by:  $\text{PMI}(i, j) = \log \frac{p(i, j)}{p(i)p(j)}$  where  $p(i, j)$  represents the co-occurrence probability of word  $i$  and  $j$  and estimated by  $p(i, j) = \frac{\#Co-occurrence}{\#Windows}$ ,  $p(i)$  represents the

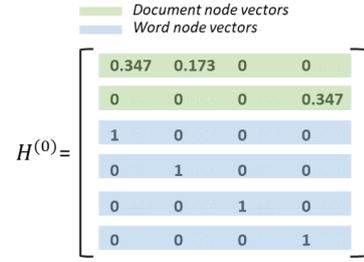


Fig. 1: Input Vectors Representations when two input documents are “word<sub>1</sub> word<sub>1</sub> word<sub>2</sub> word<sub>3</sub>” and “word<sub>3</sub> word<sub>4</sub>”.

probability of word  $i$  and estimated by  $p(i) = \frac{\#Occurrence}{\#Windows}$ . The graph is un-directed and all the edges are symmetrical.

### D. InducT-GCN Learning and Testing

After building the graph, we train it using a two-layer GCN as in [7]. The first GCN layer learns the word embeddings. The dimension of the second GCN layer is the number of classes of the dataset, and the output is fed into a softmax activation function. For example, a binary classification task will result in imension of the second GCN layer as 2. The node representations on the training documents are used for cross-entropy loss calculation and back-propagation. Formally, the propagation can be described as:

$$H^{(1)} = \sigma(\tilde{A}H^{(0)}W^{(0)}) \quad (2)$$

$$Z = \text{softmax}(\tilde{A}H^{(1)}W^{(1)}) \quad (3)$$

where  $W^{(0)}$  is a learned word embedding lookup table, and  $W^{(1)}$  is learned weight matrix in second layer. Loss is calculated by using cross-entropy between  $Z_i$  and  $Y_i, \forall i \in V_{train}$ .

In GCN, the propagation for each layer is conducted by updating the nodes with the weighted sum of their first-order neighbours and the node itself. In order to make predictions on the test set, the first-order and second-order neighbours’ representation for each test document should be aggregated. Note that we utilise the test documents during the testing phase only, so there is no need to update all the nodes in the graph during the propagation.

Instead, we conduct an one-direction propagation and only update test document nodes. Firstly,  $H_i^{(1)}, \forall i \in V_{word}$ ,  $W^{(0)}$  and  $W^{(1)}$  are recorded after training phase, and  $H_i^{(1)}, \forall i \in V_{word}$  is notated as  $H_{word}^{(1)}$ . Storing  $H_{word}^{(1)}$  enables InducT-GCN not to work on the training document nodes during the test phase, so it saves computation resources. During the testing phase, InducT-GCN supports batch testing [11]. For each batch of test document node  $B \in V_{test}, |B| = b$ , the edges  $E_B$  between  $B$  and  $V_{word}$  are calculated using TF-IDF with the document frequency of the training set.

---

**Algorithm 1** InducT-GCN Training and Testing Phase

---

**Input:** Training Graph  $G(V, \tilde{A})$ ,  $V = \{V_{train}, V_{word}\}$ ;  
Training input vectors  $H^{(0)}$ ;  
Training Labels  $\{Y_i, \forall i \in V_{train}\}$ ;  
Adjacency Matrix for Test Batch Subgraph  $\{A_B, B \in V_{test}\}$ ;  
Test input vectors  $\{H_B^{(0)}, B \in V_{test}\}$   
**Parameter:** Weight matrices  $W^{(0)}$  and  $W^{(1)}$   
**Output:** Prediction Results  $\{Y_B, \forall B \in V_{test}\}$

- 1: **for**  $epoch = 1, 2, \dots$  **do**
- 2:    $H^{(1)} \leftarrow \sigma(\tilde{A}H^{(0)}W^{(0)})$
- 3:    $Z \leftarrow softmax(\tilde{A}H^{(1)}W^{(1)})$
- 4:    $L \leftarrow CrossEntropy(Y_i, Z_i), \forall i \in V_{train}$
- 5:   Update  $W^{(0)}$  and  $W^{(1)}$
- 6: **end for**
- 7:  $H_{word}^{(1)} \leftarrow H_i^{(1)}, \forall i \in V_{word}$
- 8: **for** Batch  $B$  in  $V_{test}$  **do**
- 9:    $G' \leftarrow \{A_B, H_{word}^{(0)}, H_B^{(0)}\}$
- 10:    $H_B^{(1)} \leftarrow GCN(G', W^{(0)})$
- 11:    $G'' \leftarrow \{A_B, H_{word}^{(1)}, H_B^{(1)}\}$
- 12:    $Y_B \leftarrow argmax(GCN(G'', W^{(1)}))$
- 13: **end for**

---

Test document input  $H_B^{(0)}$  is also calculated using TF-IDF. The testing phase can be described as:

$$A_B = concat(E_B, I) \quad (4)$$

$$H_{word,B}^{(0)} = concat(H_{word}^{(0)}, H_B^{(0)}) \quad (5)$$

$$H_B^{(1)} = \sigma(A_B H_{word,B}^{(0)} W^{(0)}) \quad (6)$$

$$H_{word,B}^{(1)} = concat(H_{word}^{(1)}, H_B^{(1)}) \quad (7)$$

$$Z_B = softmax(A_B H_{word,B}^{(1)} W^{(1)}) \quad (8)$$

$$Y_B = argmax(Z_B) \quad (9)$$

where  $A_B \in R^{b \times (|V_{word}|+b)}$  stands for the weights of the weighted sum calculation and it can be considered as an adjacency matrix for test batch subgraph.  $H_B^{(0)} \in R^{(|V_{word}|+b) \times |V_{word}|}$  stands for the test batch input vectors in the subgraph.  $H_B^{(1)} \in R^{b \times h}$  is the updated test document node embedding after the first layer of GCN, and  $h$  is the hidden dimension size. Then, the stored  $H_{word}^{(1)}$  on the word nodes, which contains the first layer training documents information, are used to propagate training documents information to the test document nodes in the second layer. We formally describe the overall algorithms for the training and testing phase of InducT-GCN in Algorithm 1.

### E. Space and Time Analysis

Compared with TextGCN, InducT-GCN is more efficient both in time and space. For the space complexity: (1) Number of Parameters of InducT-GCN is  $|V_{word}| * h + h * c$  while TextGCN requires  $(|V_{train}| + |V_{word}| + |V_{test}|) * h + h * c$  parameters. Meanwhile,  $|V_{word}|$  in InducT-GCN is smaller than that in TextGCN. (2) Graph Space complexity of InducT-GCN is  $O(|V_{word}|^2 + |V_{train}| * |V_{word}|)$  and for TextGCN,

Dataset	# Train	# Test	# Word	# Class	Avg Len
R8	274	2,189	1,878	8	62.22
R52	326	2,568	2,568	52	66.98
Ohsumed	167	4,043	2,667	23	123.7
20NG	113	7,532	2,839	20	163.5
MR	355	3,554	605	2	7.25

TABLE I: Summary statistics of datasets.

it is  $O(|V_{word}|^2 + (|V_{train}| + |V_{test}|) * |V_{word}|)$ . Similarly,  $|V_{word}|$  in InducT-GCN is smaller than TextGCN. Compared with TextGCN, our InducT-GCN is faster in three ways: (1) When constructing the graph, the time complexity of PMI is  $O(|V_{word}|^2 * \#Windows)$ , and it is smaller for InducT-GCN. (2) The training time is shorter for InducT-GCN since the graph is smaller. (3) When testing on new samples, TextGCN requires retraining while InducT-GCN can make predictions without retraining.

## IV. EVALUATION SETUP

We evaluate the performance of our InducT-GCN on text classification and examine the effectiveness of the proposed inductive learning approach.

### A. Dataset

We first evaluate InducT-GCN on 5(five) publicly available text classification benchmark datasets, including R8, R52, Ohsumed, 20NG, and MR. To test in the limited labelled data environment, we select 5% of the full training set (1% for 20NG, due to the size) and remain the original test size. The detailed statistics of datasets can be found in Table I. We also apply the data augmentation methods on R8 test set to evaluate on the larger test sets, called R8A. **R8**, **R52** are two subsets of the Reuters dataset and focus on the topic classification. **Ohsumed** is produced by the MEDLINE database, containing cardiovascular diseases abstracts. **20NG**(20 NewsGroup) is a 20 class-based news classification dataset. **MR**(Movie Review) is a binary (positive and negative) sentiment polarity analysis. **R8A**: To evaluate our InducT-GCN scalability in the larger test set, we apply a data augmentation technique. Nlpaug [32] is applied for augmenting the R8 test set. To achieve this, we randomly choose one of the four options: (1) randomly deleting a word, (2) adding a word based on Word2Vec embedding similarity, (3) substituting a word using Synonyms in WordNet [33], (4) randomly swapping two words. The detailed evaluation can be found in Section V-B.

All datasets are preprocessed based on [1]. We remove the words if shown less than twice in the training documents since words only shown once can not work as a bridge between two document nodes. Words listed in the NLTK stopwords list are also removed. We apply the same preprocessing method for all experiments.

### B. Baselines

We compare InducT-GCN with baselines, mainly those models with no external resources and learning inductively. However, due to the limited number of baselines, we include

Method	PT.	R8 5%	R52 5%	Ohsumed 5%	20NG 1%	MR 5%
TF-IDF + SVM	✗	0.8054 ± 0.0000	0.6830 ± 0.0000	0.1476 ± 0.0000	0.1289 ± 0.0000	0.5537 ± 0.0000
TFIDF + LR	✗	0.8090 ± 0.0000	0.6869 ± 0.0000	0.1813 ± 0.0000	0.1860 ± 0.0000	0.5967 ± 0.0000
CNN-rand	✗	0.8107 ± 0.0110	0.6854 ± 0.0100	0.1586 ± 0.0079	0.1390 ± 0.0179	0.5485 ± 0.0122
CNN (Pretrain)	✓	0.9052 ± 0.0097	0.7708 ± 0.0181	0.3411 ± 0.0370	0.2969 ± 0.0277	<b>0.7009 ± 0.0060</b>
LSTM-rand	✗	0.7392 ± 0.0146	0.6364 ± 0.0060	0.1614 ± 0.0085	0.0766 ± 0.0063	0.5301 ± 0.0191
LSTM (Pretrain)	✓	0.7916 ± 0.0499	0.6667 ± 0.0303	0.2486 ± 0.0392	0.1010 ± 0.0220	0.6680 ± 0.0198
TextGCN [6]	✗	0.9116 ± 0.0127	0.7885 ± 0.0751	0.2225 ± 0.1138	0.2198 ± 0.1293	0.5341 ± 0.0216
SGC [8]	✗	0.8955 ± 0.0098	0.7725 ± 0.0189	0.2474 ± 0.0392	0.2948 ± 0.0342	0.6015 ± 0.0051
TextING [25]	✓	0.8648 ± 0.0414	0.7465 ± 0.0298	0.3026 ± 0.0235	N/A	0.6117 ± 0.0342
InducT-SGC	✗	0.9045 ± 0.0046	0.8046 ± 0.0066	0.3106 ± 0.0061	0.2990 ± 0.0251	0.6017 ± 0.0048
InducT-GCN	✗	<b>0.9155 ± 0.0051</b>	<b>0.8135 ± 0.0384</b>	<b>0.3563 ± 0.0078</b>	<b>0.3461 ± 0.0337</b>	0.6037 ± 0.0038

TABLE II: Comparison of with Baseline on Limited Labeled Data. For 20NG dataset, TextING [25] has out of Memory issue and they also have not tested on the 20NG either. \*The column PT. refers to the model applied any pre-trained embedding.

the baselines with pre-trained word embeddings, such as CNN (Pretrain), LSTM (Pretrain), and TextING. We also add transductive models, including TextGCN and SGC.

- **TF-IDF + SVM/LR** applies TF-IDF vectors and uses Support Vector Machine (SVM) or Logistic Regression (LR) as classifiers respectively.
- **CNN/LSTM** [1], [34] apply CNN and Long Short-Term Memory with randomly initialised word embeddings or pretrained GloVe [27] embeddings.
- **TextGCN/SGC/TextING** [6], [8], [25] are GNN text classification models. TextGCN and SGC are corpus-level GNN and TextING is document-level GNN.

### C. Settings

We apply the same set of hyper-parameters to all datasets without hyper-parameter tuning for a fair comparison. For TextGCN [6], SGC [8], our InducT-GCN and InducT-SGC, as described in [6], we applied two layers graph convolutional, and the hidden dimension is 200. Adam optimizer with a learning rate of 0.02 is used for training. For each experiment, followed by [7], 200 epochs are set to be the maximum number of epochs, and early stopping of 10 epochs is applied. 10% of the training set is randomly selected as the validation set. An early stopping mechanism is also applied for other baseline models by using the default hyperparameters.

Followed by [6], [8], [25], we use the accuracy as an evaluation metric and produce the average and standard deviation of the ten-time running results for each testing result.

## V. RESULT

### A. Performance Evaluation

A comprehensive experiment is conducted on the five benchmark datasets in the limited environment as mentioned in Section IV-A. The result presented in Table II shows that our proposed InducT-GCN significantly outperforms all baselines in terms of average accuracy on four datasets in R8, R52, Ohsumed, 20NG. Note that the baselines include transductive graph-based models, such as TextGCN and SGC, and CNN, LSTM, TextING use external resources, like pre-trained word embeddings. Meanwhile, the standard derivation is smaller

than most baseline models, showing the robustness of our model.

For more in-depth performance analysis comparing our InducT-GCN with baseline models, we can highlight that this result shows the effectiveness of the proposed InducT-GCN on long text datasets. While the average lengths of 20NG and Ohsumed are longer than 100 and those of R8 and R52 are still longer than 60, MR is less than 10, which can be considered as an extremely short text dataset. We found that models with pre-trained word embeddings GloVe perform better on short text documents. This is mainly because it would be challenging to recognise the global word co-occurrence with this short text document length, leading to fewer connections (bridging) between word nodes in corpus-level text graphs. Nevertheless, our InducT-GCN performs the best among the models with no pre-trained embeddings.

With our Inductive graph construction and learning framework, it is possible to expand to any corpus-level and transductive GCN-based text classification models, such as TextGCN [6], SGC [8], TensorGCN [9], and S<sup>2</sup>GC [10]. This section reports the generalisation capability of our inductive graph construction and learning framework. We now expand our inductive framework to another corpus-level graph-based model, SGC [8], and called InducT-SGC. Table II also visualises the comparison of the original transductive SGC models and our InducT-SGC. As shown in the table, our InducT-SGC produces much higher performance than the original SGC when the labeled data are limited. The performance improvement between both pairs of the original transductive and our inductive model, TextGCN-to-InducT-GCN and SGC-to-InducT-SGC, clearly shows the generalisation capability of our proposed inductive framework. It can also be applied to other corpus-level graph-based text classification models in the future.

### B. Impact of Test Size

As mentioned earlier, we use the R8A (R8 with a data augmentation) to show the scalability of our proposed Inductive learning framework by comparing TextGCN and InducT-GCN in the larger text set. Figure 2a shows the comparison of TextGCN and InducT-GCN on different test sizes, with 1

Embedding	R8 5%	R52 5%	Ohsumed 5%	20NG 1%	MR 5%
Random	0.9155 ± 0.0051	0.8135 ± 0.0384	<b>0.3563 ± 0.0078</b>	0.3461 ± 0.0337	0.6037 ± 0.0038
Pretrained Word2Vec	0.9124 ± 0.0043	<b>0.8290 ± 0.0084</b>	0.3544 ± 0.0305	0.3476 ± 0.0086	0.6003 ± 0.0045
Pretrained GloVe	<b>0.9159 ± 0.0102</b>	0.8266 ± 0.0090	0.3514 ± 0.0186	<b>0.3662 ± 0.0197</b>	<b>0.6051 ± 0.0055</b>

TABLE III: Test Accuracy with Different Initialized Embedding Method

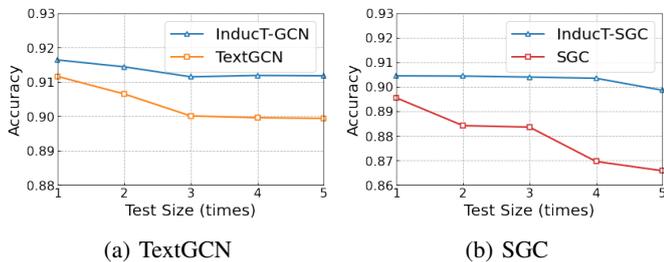


Fig. 2: Test accuracy with different test size on R8A by using TextGCN, SGC and our proposed Inductive model.

to 5 times (2,189, 4,378, 6,567, 8,756, 10,945) of the R8 original test set size. The larger the test size, the larger the gap between TextGCN and InducT-GCN. TextGCN produces worse performance with the largest test size. This is mainly because only a small proportion of the document nodes would contribute to the gradient in TextGCN with a larger test set. Especially during the training phase, it is difficult for TextGCN to learn embeddings of those test document nodes having fewer connections with word nodes by backpropagation. Moreover, we found that the performance of our InducT-GCN decreased only a little (less than 0.5) when the test size grew. We also conducted the same evaluation based on SGC by applying our Inductive graph construction and learning framework to SGC, InducT-SGC. Like the result that our InducT-GCN produced, the InducT-SGC delivers a much higher performance than the original SGC. The performance trend shows how our Inductive framework perfectly fits the inductive learning nature.

### C. Impact of Initial Word Embedding

As mentioned in Section III-D, the first layer of InducT-GCN learns the word embeddings and is randomly initialized. We also examine other initial embedding weights methods including pre-trained Word2vec [26] and GloVe [27]. Table III shows the performance comparison of different pre-trained word embeddings. In most cases except Ohsumed, pre-trained word embeddings can improve the performance of InducT-GCN. However, Ohsumed is a medical-related dataset, and out-of-vocabulary issue of the pretrained embedding doesn't help on the document classification task. Although we only focus on the model without using any external resources in this study, this result still shows the potentiality of InducT-GCN when used with external resources.

### D. Computation Time Results

Table IV compares the original transductive TextGCN with InducT-GCN on R8A and visualises the superiority of our inductive learning framework by reducing the computation

Test Size	TextGCN		InducT-GCN	
	Graph	Training	Graph	Training
×1	6.29	2.75	0.89	0.52
×2	11.90	3.20	1.11	0.53
×3	16.60	3.54	1.30	0.53
×4	21.10	4.13	1.52	0.55
×5	27.50	4.98	1.68	0.51

TABLE IV: Graph Construction Time and Training Time comparison on R8A (sec). Hardware: 16 Intel(R) Core(TM) i9-9900X CPU @ 3.50GHz and NVIDIA Titan RTX 24GB

Method	R8 Full	R52 Full
TextGCN	0.9629 ± 0.0010	0.9295 ± 0.0020
InducT-GCN	<b>0.9653 ± 0.0017</b>	<b>0.9323 ± 0.0015</b>

TABLE V: Test Accuracy on Full Data Setting

time, including graph construction and training. The larger the test size is, the more time InducT-GCN can save.

### E. Performance in Full Dataset

We also evaluated the performance of our InducT-GCN with the full dataset, like the TextGCN was evaluated [6]. As can be seen in Table V, the performance of InducT-GCN and TextGCN on R8 and R52 are comparable when using the entire dataset with the same hyperparameters. We can conclude that InducT-GCN is superior to the TextGCN in terms of performance and computation, which is not only in smaller space with fewer parameter numbers but also in the whole dataset setting.

## VI. CONCLUSION

This study proposes a novel inductive graph-based text classification framework, InducT-GCN, which makes heavy and transductive GCN-based text classification models inductive. We construct a graph only using training set statistics. InducT-GCN can efficiently capture global information with fewer parameters and smaller space complexity. Our InducT-GCN significantly outperformed all graph-based text classification baselines and was even better than the models using pretrained embeddings. We also demonstrated the generalisation capability of our inductive graph construction and learning framework by applying and expanding different transductive graph-based text classification models, like TextGCN and SGC. Compared to the original models, the performance and computation time were surprisingly improved. It is hoped that this paper provides some insight into the future integration of the lighter and faster inductive graph neural networks on different NLP tasks.

## REFERENCES

- [1] Y. Kim, "Convolutional neural networks for sentence classification," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014*, 2014, pp. 1746–1751.
- [2] P. Zhou, Z. Qi, S. Zheng, J. Xu, H. Bao, and B. Xu, "Text classification improved by integrating bidirectional lstm with two-dimensional max pooling," in *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, 2016, pp. 3485–3495.
- [3] H. Peng, J. Li, Y. He, Y. Liu, M. Bao, L. Wang, Y. Song, and Q. Yang, "Large-scale hierarchical text classification with recursively regularized deep graph-cnn," in *Proceedings of the 2018 world wide web conference*, 2018, pp. 1063–1072.
- [4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [5] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized BERT pretraining approach," *CoRR*, vol. abs/1907.11692, 2019.
- [6] L. Yao, C. Mao, and Y. Luo, "Graph convolutional networks for text classification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 7370–7377.
- [7] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- [8] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger, "Simplifying graph convolutional networks," in *International conference on machine learning*. PMLR, 2019, pp. 6861–6871.
- [9] X. Liu, X. You, X. Zhang, J. Wu, and P. Lv, "Tensor graph convolutional networks for text classification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 05, 2020, pp. 8409–8416.
- [10] H. Zhu and P. Koniusz, "Simple spectral graph convolution," in *International Conference on Learning Representations*, 2021.
- [11] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 1025–1035.
- [12] J. Chen, T. Ma, and C. Xiao, "Fastgcn: Fast learning with graph convolutional networks via importance sampling," in *International Conference on Learning Representations*, 2018.
- [13] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph Attention Networks," *International Conference on Learning Representations*, 2018.
- [14] R. A. Rossi, R. Zhou, and N. K. Ahmed, "Deep inductive graph representation learning," *IEEE Computer Architecture Letters*, vol. 32, no. 03, pp. 438–452, 2020.
- [15] da Xu, chuanwei ruan, evren korpeoglu, sushant kumar, and kannan achan, "Inductive representation learning on temporal graphs," in *International Conference on Learning Representations (ICLR)*, 2020.
- [16] J. Bastings, I. Titov, W. Aziz, D. Marcheggiani, and K. Sima'an, "Graph convolutional encoders for syntax-aware neural machine translation," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017, pp. 1957–1967.
- [17] M. Tu, G. Wang, J. Huang, Y. Tang, X. He, and B. Zhou, "Multi-hop reading comprehension across multiple documents by reasoning over heterogeneous graphs," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 2704–2713.
- [18] Y. Li, R. Jin, and Y. Luo, "Classifying relations in clinical narratives using segment graph convolutional and recurrent neural networks (seg-gcrns)," *Journal of the American Medical Informatics Association: JAMIA*, vol. 26, no. 3, pp. 262–268, 2019.
- [19] Y. Cao, Z. Liu, C. Li, J. Li, and T.-S. Chua, "Multi-channel graph neural network for entity alignment," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 1452–1461.
- [20] T. Yang, L. Hu, C. Shi, H. Ji, X. Li, and L. Nie, "Hgat: Heterogeneous graph attention networks for semi-supervised short text classification," *ACM Transactions on Information Systems*, vol. 39, no. 3, May 2021.
- [21] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *NIPS*, 2016.
- [22] Y. Zhang, Q. Liu, and L. Song, "Sentence-state lstm for text representation," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 317–327.
- [23] G. Nikolentzos, A. Tixier, and M. Vazirgiannis, "Message passing attention networks for document understanding," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 05, 2020, pp. 8544–8551.
- [24] L. Huang, D. Ma, S. Li, X. Zhang, and W. Houfeng, "Text level graph neural network for text classification," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 3435–3441.
- [25] Y. Zhang, X. Yu, Z. Cui, S. Wu, Z. Wen, and L. Wang, "Every document owns its structure: Inductive text classification via graph neural networks," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 334–339.
- [26] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2013.
- [27] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [28] C. Li, X. Peng, H. Peng, J. Li, and L. Wang, "Textgl: Graph-based transductive learning for semi-supervised text classification via structure-sensitive interpolation," in *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI)*, 2021, pp. 2680–2686.
- [29] H. Zhang and J. Zhang, "Text graph transformer for document classification," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 8322–8327.
- [30] Y. Lin, Y. Meng, X. Sun, Q. Han, K. Kuang, J. Li, and F. Wu, "Bertgcn: Transductive text classification by combining gcn and bert," *arXiv preprint arXiv:2105.05727*, 2021.
- [31] S. Aji and R. Kaimal, "Document summarization using positive pointwise mutual information," *AIRCC's International Journal of Computer Science and Information Technology*, vol. 4, no. 2, pp. 47–55, 2012.
- [32] E. Ma, "Nlp augmentation," <https://github.com/makcedward/nlpaug>, 2019.
- [33] G. A. Miller, "Wordnet: a lexical database for english," *Communications of the ACM*, vol. 38, no. 11, pp. 39–41, 1995.
- [34] P. Liu, X. Qiu, and X. Huang, "Recurrent neural network for text classification with multi-task learning," in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*, S. Kambhampati, Ed. IJCAI/AAAI Press, 2016, pp. 2873–2879.