

# Be More with Less: Hypergraph Attention Networks for Inductive Text Classification

**Kaize Ding**

Arizona State University  
kaize.ding@asu.edu

**Jianling Wang**

Texas A&M University  
jllwang@tamu.edu

**Jundong Li**

University of Virginia  
jundong@virginia.edu

**Dingcheng Li**

Amazon Inc.  
lidingch@amazon.com

**Huan Liu**

Arizona State University  
huan.liu@asu.edu

## Abstract

Text classification is a critical research topic with broad applications in natural language processing. Recently, graph neural networks (GNNs) have received increasing attention in the research community and demonstrated their promising results on this canonical task. Despite the success, their performance could be largely jeopardized in practice since they are: (1) unable to capture high-order interaction between words; (2) inefficient to handle large datasets and new documents. To address those issues, in this paper, we propose a principled model – hypergraph attention networks (HyperGAT), which can obtain more expressive power with less computational consumption for text representation learning. Extensive experiments on various benchmark datasets demonstrate the efficacy of the proposed approach on the text classification task.

## 1 Introduction

Text classification, as one of the most fundamental tasks in the field of natural language processing, has received continuous endeavors from researchers due to its wide spectrum of applications, including sentiment analysis (Wang et al., 2016), topic labeling (Wang and Manning, 2012), and disease diagnosis (Miotto et al., 2016). Inspired by the success of deep learning techniques, methods based on representation learning such as convolutional neural networks (CNNs) (Kim, 2014) and recurrent neural networks (RNNs) (Liu et al., 2016) have been extensively explored in the past few years. In essence, the groundbreaking achievements of those methods can be attributed to their strong capability of capturing sequential context information from local consecutive word sequences.

More recently, graph neural networks (GNNs) (Kipf and Welling, 2017; Veličković et al., 2018; Hamilton et al., 2017) have drawn

much attention and demonstrated their superior performance in the task of text classification (Yao et al., 2019; Wu et al., 2019a; Liu et al., 2020). This line of work leverages the knowledge from both training and test documents to construct a corpus-level graph with global word co-occurrence and document-word relations, and consider text classification as a semi-supervised node classification problem. Then with GNNs, long-distance interactions between words could be captured to improve the final text classification performance.

Despite their promising early results, the usability of existing efforts could be largely jeopardized in real-world scenarios, mainly owing to their limitations in the following two aspects: (i) *Expressive Power*. Existing GNN-based methods predominantly focus on pairwise interactions (i.e., dyadic relations) between words. However, word interactions are not necessarily dyadic in natural language, but rather could be triadic, tetradic, or of a higher-order. For instance, consider the idiom “eat humble pie”, whose definition is “admit that one was wrong” in common usage. If we adopt a simple graph to model the word interactions, GNNs may misinterpret the word `pie` as “a baked dish” based on its pairwise connections to other two words (`humble-pie` and `eat-pie`), then further misunderstand the actual meaning of the whole idiom. Hence, how to go beyond pairwise relations and further capture the high-order word interactions is vital for high-quality text representation learning, but still remains to be explored. (ii) *Computational Consumption*. On the one hand, most of the endeavors with GNN backbone tend to be memory-inefficient when the scale of data increases, due to the fact that constructing and learning on a global document-word graph consumes immense memory (Huang et al., 2019). On the other hand, the mandatory access to test documents during training renders those methods inherently *trans-*

*ductive*. It means that when new data arrives, we have to retrain the model from scratch for handling newly added documents. Therefore, it is necessary to design a computationally efficient approach for solving graph-based text classification.

Upon the discussions above, one critical research question to ask is “*Is it feasible to acquire more expressive power with less computational consumption?*”. To achieve this goal, we propose to adopt document-level hypergraph (hypergraph is a generalization of simple graph, in which a hyperedge can connect *arbitrary number of nodes*) for modeling each text document. The use of document-level hypergraphs potentially enables a learning model not only to alleviate the computational inefficiency issue, but more remarkably, to capture heterogeneous (e.g., sequential and semantic) high-order contextual information of each word. Therefore, more expressive power could be obtained with less computational consumption during the text representation learning process. As conventional GNN models are infeasible to be used on hypergraphs, to bridge this gap, we propose a new model named HyperGAT, which is able to capture the encoded high-order word interactions within each hypergraph. In the meantime, its internal dual attention mechanism highlights key contextual information for learning highly expressive text representations. To summarize, our contributions are in three-fold:

- We propose to model text documents with document-level hypergraphs, which improves the model expressive power and reduces computational consumption.
- A principled model HyperGAT based on a dual attention mechanism is proposed to support representation learning on text hypergraphs.
- We conduct extensive experiments on multiple benchmark datasets to illustrate the superiority of HyperGAT over other state-of-the-art methods on the text classification task.

## 2 Related Work

### 2.1 Graph Neural Networks

Graph neural networks (GNNs) – a family of neural models for learning latent node representations in a graph, have achieved remarkable success in different graph learning tasks (Defferrard et al., 2016; Kipf and Welling, 2017; Veličković et al., 2016; Ding et al., 2019a, 2020). Most of the prevailing GNN models follow the paradigm of neighborhood

aggregation, aiming to learn latent node representations via message passing among local neighbors in the graph. With deep roots in graph spectral theory, the learning process of graph convolutional networks (GCNs) (Kipf and Welling, 2017) can be considered as a mean-pooling neighborhood aggregation. Later on, GraphSAGE (Hamilton et al., 2017) was developed to concatenate the node’s feature with mean/max/LSTM pooled neighborhood information, which enables inductive representation learning on large graphs. Graph attention networks (GATs) (Veličković et al., 2016) incorporate trainable attention weights to specify fine-grained weights on neighbors when aggregating neighborhood information of a node. Recent research further extend GNN models to consider global graph information (Battaglia et al., 2018) and edge information (Gilmer et al., 2017) during aggregation. More recently, hypergraph neural networks (Feng et al., 2019; Bai et al., 2020; Wang et al., 2020) are proposed to capture high-order dependency between nodes. Our model HyperGAT is the first attempt to shift the power of hypergraph to the canonical text classification task.

### 2.2 Deep Text Classification

Grounded on the fast development of deep learning techniques, various neural models that automatically represent texts as embeddings have been developed for text classification. Two representative deep neural models, CNNs (Kim, 2014; Zhang et al., 2015) and RNNs (Tai et al., 2015; Liu et al., 2016) have shown their superior power in the text classification task. To further improve the model expressiveness, a series of attentional models have been developed, including hierarchical attention networks (Yang et al., 2016), attention over attention (Cui et al., 2017), etc. More recently, graph neural networks have shown to be a powerful tool for solving the problem of text classification by considering the long-distance dependency between words. Specifically, TextGCN (Yao et al., 2019) applies the graph convolutional networks (GCNs) (Kipf and Welling, 2017) on a single large graph built from the whole corpus, which achieves state-of-the-art performance on text classification. Later on, SGC (Wu et al., 2019a) is proposed to reduce the unnecessary complexity and redundant computation of GCNs, and shows competitive results with superior time efficiency. TensorGCN (Liu et al., 2020) proposes a text graph

tensor to learn word and document embeddings by incorporating more context information. (Huang et al., 2019) propose to learn text representations on document-level graphs. However, those transductive methods are computationally inefficient and cannot capture the high-order interactions between words for improving model expressive power.

### 3 Methodology

In this section, we introduce a new family of GNN models developed for inductive text classification. By reviewing the existing GNN-based endeavors, we first summarize their main limitations that need to be addressed. Then we illustrate how we use hypergraphs to model text documents for achieving the goals. Finally, we propose the model HyperGAT based on a new dual attention mechanism and model training for inductive text classification.

#### 3.1 GNNs for Text Classification

With the booming development of deep learning techniques, graph neural networks (GNNs) have achieved great success in representation learning on graph-structured data (Zhou et al., 2018; Ding et al., 2019b). In general, most of the prevailing GNN models follow the neighborhood aggregation strategy, and a GNN layer can be defined as:

$$\mathbf{h}_i^l = \text{AGGR}^l\left(\mathbf{h}_i^{l-1}, \{\mathbf{h}_j^{l-1} | \forall j \in \mathcal{N}_i\}\right), \quad (1)$$

where  $\mathbf{h}_i^l$  is the node representation of node  $i$  at layer  $l$  (we use  $\mathbf{x}_i$  as  $\mathbf{h}_i^0$ ) and  $\mathcal{N}_i$  is the local neighbor set of node  $i$ . AGGR is the aggregation function of GNNs and has a series of possible implementations (Kipf and Welling, 2017; Hamilton et al., 2017; Veličković et al., 2018).

Given the capability of capturing long-distance interactions between entities, GNNs also have demonstrated promising performance on text classification (Yao et al., 2019; Wu et al., 2019b; Liu et al., 2020). The prevailing approach is to build a corpus-level document-word graph and try to classify documents through semi-supervised node classification. Despite their success, most of the existing efforts suffer from the computational inefficiency issue, not only because of the mandatory access of test documents, but also the construction of corpus-level document-word graphs. In the meantime, those methods are largely limited by the expressibility of using simple graphs to model word interactions. Therefore, how to improve model ex-

pressive power with less computational consumption is a challenging and imperative task to solve.

#### 3.2 Documents as Text Hypergraphs

To address the aforementioned challenges, in this study, we alternatively propose to model text documents with document-level hypergraphs. Formally, hypergraphs can be defined as follows:

**Definition 3.1 Hypergraphs:** A hypergraph is defined as a graph  $G = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V} = \{v_1, \dots, v_n\}$  represents the set of nodes in the graph, and  $\mathcal{E} = \{e_1, \dots, e_m\}$  represents the set of hyperedges. Note that for any hyperedge  $e$ , it can connect two or more nodes (i.e.,  $\sigma(e) \geq 2$ ).

Notably, the topological structure of a hypergraph  $G$  can also be represented by an incidence matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$ , with entries defined as:

$$\mathbf{A}_{ij} = \begin{cases} 1, & \text{if } v_i \in e_j, \\ 0, & \text{if } v_i \notin e_j. \end{cases} \quad (2)$$

In the general case, each node in hypergraphs could come with a  $d$ -dimensional attribute vector. Therefore, all the node attributes can be denoted as  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ , and we can further use  $G = (\mathbf{A}, \mathbf{X})$  to represent the whole hypergraph for simplicity.

For a text hypergraph, nodes represent words in the document and node attributes could be either one-hot vector or the pre-trained word embeddings (e.g., word2vec, GloVe). In order to model heterogeneous high-order context information within each document, we include multi-relational hyperedges as follows:

**Sequential Hyperedges.** Sequential context depicts the language property of local co-occurrence between words, which has demonstrated its effectiveness for text representation learning (Yao et al., 2019). To leverage the sequential context information of each word, we first construct sequential hyperedges for each document in the corpus. One natural way is to adopt a fixed-size sliding window to obtain global word co-occurrence as the sequential context. Inspired by the success of hierarchical attention networks (Yang et al., 2016), here we consider each sentence as a hyperedge and it connects all the words in this sentence. As another benefit, using sentences as sequential hyperedges enables our model to capture the document structural information at the same time.

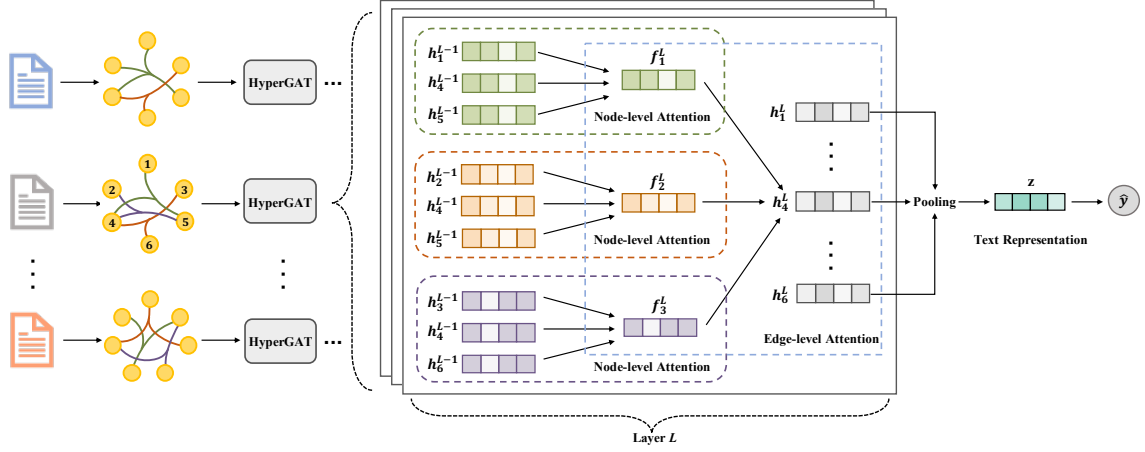


Figure 1: Illustration of the proposed hypergraph attention networks (HyperGAT) for inductive text classification. We construct a hypergraph for each text document and feed it into HyperGAT. Based on the node and edge-level attention, text representations that capture high-order word interactions can be derived. Figure best viewed in color.

**Semantic Hyperedges.** Furthermore, in order to enrich the semantic context for each word, we build semantic hyperedges to capture topic-related high-order correlations between words (Linmei et al., 2019). Specifically, we first mine the latent topics  $T$  from the text documents using LDA (Blei et al., 2003) and each topic  $t_i = (\theta_1, \dots, \theta_w)$  ( $w$  denotes the vocabulary size) can be represented by a probability distribution over the words. Then for each topic, we consider it as a semantic hyperedge that connects the top  $K$  words with the largest probabilities in the document. With those topic-related hyperedges, we are able to enrich the high-order semantic context of words in each document.

It is worth mentioning that though we only discuss sequential and semantic hyperedges in this study, other meaningful hyperedges (e.g., syntactic-related) could also be integrated into the proposed model for further improving the model expressiveness and we leave this for future work.

### 3.3 Hypergraph Attention Networks

To support text representation learning on the constructed text hypergraphs, we then propose a new model called HyperGAT (as shown in Figure 1) in this section. Apart from conventional GNN models, HyperGAT learns node representations with two different aggregation functions, allowing to capture heterogeneous high-order context information of words on text hypergraphs. In general, a HyperGAT layer can be defined as:

$$\begin{aligned} \mathbf{h}_i^l &= \text{AGGR}_{edge}^l \left( \mathbf{h}_i^{l-1}, \{ \mathbf{f}_j^l | \forall e_j \in \mathcal{E}_i \} \right), \\ \mathbf{f}_j^l &= \text{AGGR}_{node}^l \left( \{ \mathbf{h}_k^{l-1} | \forall v_k \in e_j \} \right), \end{aligned} \quad (3)$$

where  $\mathcal{E}_i$  denotes the set of hyperedges connected to node  $v_i$  and  $\mathbf{f}_j^l$  is the representation of hyperedge  $e_j$  in layer  $l$ .  $\text{AGGR}_{edge}$  is an aggregation function that aggregates features of hyperedges to nodes and  $\text{AGGR}_{node}$  is another aggregation function that aggregates features of nodes to hyperedges. In this work, we propose to implement those two functions based on a dual attention mechanism. We will start by describing a single layer  $l$  for building arbitrary HyperGAT architectures as follows:

**Node-level Attention.** Given a specific node  $v_i$ , our HyperGAT layer first learns the representations of all its connected hyperedges  $\mathcal{E}_i$ . As not all the nodes in a hyperedge  $e_j \in \mathcal{E}_i$  contribute equally to the hyperedge meaning, we introduce attention mechanism (i.e., node-level attention) to highlight those nodes that are important to the meaning of the hyperedge and then aggregate them to compute the hyperedge representation  $\mathbf{f}_j^l$ . Formally:

$$\mathbf{f}_j^l = \sigma \left( \sum_{v_k \in e_j} \alpha_{jk} \mathbf{W}_1 \mathbf{h}_k^{l-1} \right), \quad (4)$$

where  $\sigma$  is the nonlinearity such as ReLU and  $\mathbf{W}_1$  is a trainable weight matrix.  $\alpha_{jk}$  denotes the attention coefficient of node  $v_k$  in the hyperedge  $e_j$ , which can be computed by:

$$\begin{aligned} \alpha_{jk} &= \frac{\exp(\mathbf{a}_1^T \mathbf{u}_k)}{\sum_{v_p \in e_j} \exp(\mathbf{a}_1^T \mathbf{u}_p)}, \\ \mathbf{u}_k &= \text{LeakyReLU}(\mathbf{W}_1 \mathbf{h}_k^{l-1}), \end{aligned} \quad (5)$$

where  $\mathbf{a}_1^T$  is a weight vector (a.k.a, context vector).

**Edge-level Attention.** With all the hyperedges representations  $\{\mathbf{f}_j^l | \forall e_j \in \mathcal{E}_i\}$ , we again apply an edge-level attention mechanism to highlight the informative hyperedges for learning the next-layer representation of node  $v_i$ . This process can be formally expressed as:

$$\mathbf{h}_i^l = \sigma \left( \sum_{e_j \in \mathcal{E}_i} \beta_{ij} \mathbf{W}_2 \mathbf{f}_j^l \right), \quad (6)$$

where  $\mathbf{h}_i^l$  is the output representation of node  $v_i$  and  $\mathbf{W}_2$  is a weight matrix.  $\beta_{ij}$  denotes the attention coefficient of hyperedge  $e_j$  on node  $v_i$ , which can be computed by:

$$\beta_{ij} = \frac{\exp(\mathbf{a}_2^T \mathbf{v}_j)}{\sum_{e_p \in \mathcal{E}_i} \exp(\mathbf{a}_2^T \mathbf{v}_p)}, \quad (7)$$

$$\mathbf{v}_j = \text{LeakyReLU}([\mathbf{W}_2 \mathbf{f}_j^l || \mathbf{W}_1 \mathbf{h}_i^{l-1}]),$$

where  $\mathbf{a}_2^T$  is another weight (context) vector for measuring the importance of the hyperedges and  $||$  is the concatenation operation.

The proposed dual attention mechanism enables a HyperGAT layer not only to capture the high-order word interactions, but also to highlight the key information at different granularities during the node representation learning process.

### 3.4 Inductive Text Classification

For each document, after going through  $L$  HyperGAT layers, we are able to compute all the node representations on the constructed text hypergraph. Then we apply the *mean-pooling* operation on the learned node representations  $\mathbf{H}^L$  to obtain the document representation  $\mathbf{z}$ , and feed it to a softmax layer for text classification. Formally:

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{W}_c \mathbf{z} + \mathbf{b}_c), \quad (8)$$

where  $\mathbf{W}_c$  is a parameter matrix mapping the document representation into an output space and  $\mathbf{b}_c$  is the bias.  $\hat{\mathbf{y}}$  denotes the predicted label scores. Specifically, the loss function of text classification is defined as the cross-entropy loss:

$$\mathcal{L} = - \sum_d \log(\hat{\mathbf{y}}_j^d), \quad (9)$$

where  $j$  is the ground truth label of document  $d$ . Thus HyperGAT can be learned by minimizing the above loss function over all the labeled documents.

Note that HyperGAT eliminates the mandatory access of test documents during training, making

Dataset	20NG	R8	R52	Ohsumed	MR
# Doc	18,846	7,674	9,100	7,400	10,662
# Train	11,314	5,485	6,532	3,357	7,108
# Test	7,532	2,189	2,568	4,043	3,554
# Word	42,757	7,688	8,892	14,157	18,764
Avg Len	221.26	65.72	69.82	135.82	20.39
# Class	20	8	52	23	2

Table 1: Summary statistics of the evaluation datasets.

the model different from existing GNN-based methods. For unseen documents, we can directly feed their corresponding text hypergraphs to the previously learned model and compute their labels. Hence, we can handle the newly added data in an inductive way instead of retraining the model.

## 4 Experiments

### 4.1 Experimental Setting

**Evaluation Datasets.** To conduct a fair and comprehensive evaluation, we adopt five benchmark datasets from different domains in our experiments: 20-Newsgroups (20NG), Reuters (R8 and R52), Ohsumed, and Movie Review (MR). Those datasets have been widely used for evaluating graph-based text classification performance (Yao et al., 2019; Huang et al., 2019; Liu et al., 2020). Specifically, the 20-Newsgroups dataset and two Reuters datasets are used for news classification. The Ohsumed dataset is medical literature. The Movie Review dataset is collected for binary sentiment classification. A summary statistics of the benchmark datasets is presented in table 1 and more detailed descriptions can be found in (Yao et al., 2019). For quantitative evaluation, we follow the same train/test splits and data preprocessing procedure in (Yao et al., 2019) in our experiments. In each run, we randomly sample 90% of the training samples to train the model and use the left 10% data for validation. More details can be found in Appendix A.1.

**Compared Methods.** In our experiments, the baselines compared with our model HyperGAT can be generally categorized into three classes: (i) *word embedding-based* methods that classify documents based on pre-trained word embeddings, including fastText (Joulin et al., 2016), and more advanced methods SWEM (Shen et al., 2018) and LEAM (Wang et al., 2018); (ii) *sequence-based* methods which capture text fea-

Model	20NG	R8	R52	Ohsumed	MR
CNN-rand	0.7693 ± 0.0061	0.9402 ± 0.0057	0.8537 ± 0.0047	0.4387 ± 0.0100	0.7498 ± 0.0070
CNN-non-static	0.8215 ± 0.0052	0.9571 ± 0.0052	0.8759 ± 0.0048	0.5833 ± 0.0106	0.7775 ± 0.0072
LSTM	0.6571 ± 0.0152	0.9368 ± 0.0082	0.8554 ± 0.0113	0.4114 ± 0.0117	0.7506 ± 0.0044
LSTM (pretrain)	0.7543 ± 0.0172	0.9609 ± 0.0019	0.9048 ± 0.0086	0.5110 ± 0.0150	0.7733 ± 0.0089
Bi-LSTM	0.7318 ± 0.0185	0.9631 ± 0.0033	0.9054 ± 0.0091	0.4927 ± 0.0107	0.7768 ± 0.0086
fastText	0.7938 ± 0.0030	0.9613 ± 0.0021	0.9281 ± 0.0009	0.5770 ± 0.0049	0.7514 ± 0.0020
fastText (bigrams)	0.7967 ± 0.0029	0.9474 ± 0.0011	0.9099 ± 0.0005	0.5569 ± 0.0039	0.7624 ± 0.0012
SWEM	0.8516 ± 0.0029	0.9532 ± 0.0026	0.9294 ± 0.0024	0.6312 ± 0.0055	0.7665 ± 0.0063
LEAM	0.8191 ± 0.0024	0.9331 ± 0.0024	0.9184 ± 0.0023	0.5858 ± 0.0079	0.7695 ± 0.0045
Graph-CNN	0.8142 ± 0.0032	0.9699 ± 0.0012	0.9275 ± 0.0022	0.6386 ± 0.0053	0.7722 ± 0.0027
TextGCN (transductive)	0.8643 ± 0.0009	0.9707 ± 0.0010	0.9356 ± 0.0018	0.6836 ± 0.0056	0.7674 ± 0.0020
TextGCN (inductive)	0.8331 ± 0.0026	0.9578 ± 0.0029	0.8820 ± 0.0072	0.5770 ± 0.0035	0.7480 ± 0.0025
Text-level GNN	0.8416 ± 0.0025	0.9789 ± 0.0020	0.9460 ± 0.0030	0.6940 ± 0.0060	0.7547 ± 0.0006
HyperGAT (ours)	0.8662 ± 0.0016	0.9797 ± 0.0023	0.9498 ± 0.0027	0.6990 ± 0.0034	0.7832 ± 0.0027

Table 2: Test accuracy on document classification with different models. Each model we ran 10 times and report the mean ± standard deviation. HyperGAT significantly outperforms all the baselines based on t-tests ( $p < 0.05$ ).

tures from local consecutive word sequences, including CNNs (Kim, 2014), LSTMs (Liu et al., 2016), and Bi-LSTM (Huang et al., 2015); (iii) *graph-based* methods that aim to capture interactions between words, including Graph-CNN (Deferrard et al., 2016), two versions of TextGCN (Yao et al., 2019) and Text-level GNN (Huang et al., 2019). Note that TextGCN (transductive) is the model proposed in the original paper and TextGCN (inductive) is the inductive version implemented by the same authors. Text-level GNN is a state-of-the-art baseline which performs text representation learning on document-level graphs. More details of baselines can be found in (Yao et al., 2019).

**Implementation Details.** HyperGAT is implemented by PyTorch and optimized with the Adam optimizer. We train and test the model on a 12 GB Titan Xp GPU. Specifically, our HyperGAT model consists of two layers with 300 and 100 embedding dimensions, respectively. We use one-hot vectors as the node attributes and the batch size is set to 8 for all the datasets. The optimal values of hyperparameters are selected when the model achieves the highest accuracy for the validation samples. The optimized learning rate  $\alpha$  is set to 0.0005 for MR and 0.001 for the other datasets. L2 regularization is  $10^{-6}$  and dropout rate is 0.3 for the best performance. For learning HyperGAT, we train the model for 100 epochs with early-stopping strategy. To construct the semantic hyperedges, we train an LDA model for each dataset using the training documents and select the Top-10 words from each topic. The topic number is set to the same

number of classes. For baseline models, we either show the results reported in previous research (Yao et al., 2019) or run the codes provided by the authors using the parameters described in the original papers. *More details can be found in the Appendix A.2.* Our data and source code is available at <https://github.com/kaize0409/HyperGAT>.

## 4.2 Experimental Results

**Classification Performance.** We first conduct comprehensive experiments to evaluate model performance on text classification and present the results in Table 2. Overall, our model HyperGAT outperforms all the baselines on the five evaluation datasets, which demonstrates its superior capability in text classification. In addition, we can make the following in-depth observations and analysis:

- *Graph-based* methods, especially GNN-based models are able to achieve superior performance over the other two categories of baselines on the first four datasets. This observation indicates that text classification performance can be directly improved by capturing long-distance word interactions. While for the MR dataset, *sequence-based* methods (CNNs and LSTMs) show stronger classification capability than most of the *graph-based* baselines. One potential reason is that sequential context information plays a critical role in sentiment classification, which cannot be explicitly captured by the majority of existing *graph-based* methods.
- Not surprisingly, without the additional knowledge on test documents, the performance of

Model	TextGCN (transductive)	HyperGAT
<b>20NG</b>	1,4479.36MB	180.33MB
<b>R8</b>	931.58MB	41.75MB
<b>R52</b>	1289.48MB	46.85MB
<b>Ohsumed</b>	1822.71MB	63.17MB
<b>MR</b>	3338.24MB	80.99MB

Table 3: GPU memory consumption of different methods. The batch size for HyperGAT is set to 8.

TextGCN (inductive) largely falls behind its original transductive version. Though Text-level GNN is able to achieve performance improvements by adding trainable edge weights between word, its performance is still limited by the information loss of using pairwise simple graph. In particular, our model HyperGAT achieves considerable improvements over other GNN-based models, demonstrating the importance of high-order context information for learning word representations.

**Computational Efficiency.** Table 3 presents the computational cost comparison between the most representative transductive baseline TextGCN and our approach. From the reported results, we can clearly find that HyperGAT has a significant computational advantage in terms of memory consumption. The main reason is that HyperGAT conducts text representation learning at the document-level and it only needs to store a batch of small text hypergraphs during training. On the contrary, TextGCN requires constructing a large document-word graph using both training and test documents, which inevitably consumes a great amount of memory. Another computational advantage of our model is that HyperGAT is an inductive model that can generalize to unseen documents. Thus we do not have to retrain the whole model for newly added documents like transductive methods.

**Model Sensitivity.** The model performance on 20NG and Ohsumed with different first-layer embedding dimensions is reported in Figure 2, and we omit the results on other datasets since similar results can be observed. Notably, the best performance of HyperGAT is achieved when the first-layer embedding size is set to 300. It indicates that small embedding size may render the model less expressive, while the model may encounter overfitting if the embedding size is too large. In the meantime, to evaluate the effect of the size of

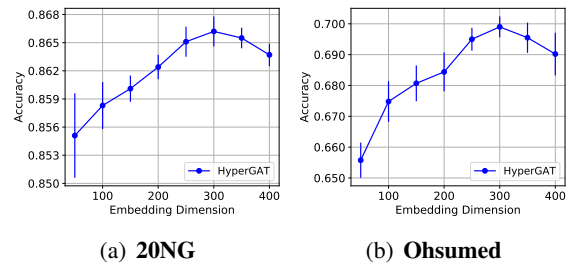


Figure 2: Test accuracy by varying the embedding size of the first HyperGAT layer.

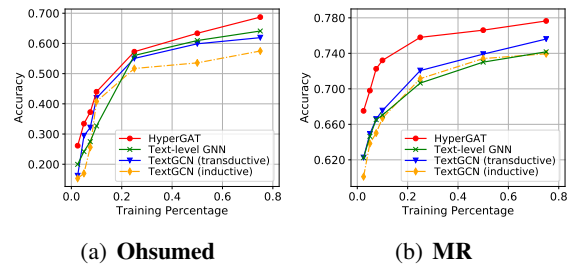


Figure 3: Test accuracy by varying the proportions of training data (2.5%, 5%, 7.5%, 10%, 25%, 50%, 75%).

labeled training data, we compare several best performing models with different proportions of the training data and report the results on Ohsumed and MR in Figure 3. In general, with the growth of labeled training data, all the evaluated methods can achieve performance improvements. More remarkably, HyperGAT can significantly outperform other baselines with limited labeled data, showing its effectiveness in real-world scenarios.

### 4.3 Ablation Analysis

To investigate the contribution of each module in HyperGAT, we conduct an ablation analysis and report the results in Table 4. Specifically, *w/o attention* is a variant of HyperGAT that replaces the dual attention with convolution. *w/o sequential* and *w/o semantic* are another two variants by excluding sequential, semantic hyperedges, respectively. From the reported results we can learn that HyperGAT can achieve better performance by stacking more layers. This observation can verify the usefulness of long-distance word interactions for text representation learning. Moreover, the performance gap between *w/o attention* and HyperGAT shows the effectiveness of the dual attention mechanism for learning more expressive word representations. By comparing the results of *w/o sequential* and *w/o semantic*, we can learn that the context informa-

Model	20NG	R8	R52	Ohsumed	MR
w/o attention	0.8645 ± 0.0006	0.9705 ± 0.0015	0.9321 ± 0.0023	0.6611 ± 0.0042	0.7699 ± 0.0044
w/o sequential	0.6813 ± 0.0024	0.9448 ± 0.0053	0.9051 ± 0.0023	0.5664 ± 0.0047	0.7766 ± 0.0009
w/o semantic	0.8602 ± 0.0031	0.9714 ± 0.0026	0.9415 ± 0.0032	0.6848 ± 0.0045	0.7811 ± 0.0028
HyperGAT (1 layer)	0.8610 ± 0.0014	0.9735 ± 0.0012	0.9472 ± 0.0023	0.6913 ± 0.0023	0.7788 ± 0.0016
HyperGAT	0.8662 ± 0.0016	0.9797 ± 0.0023	0.9498 ± 0.0027	0.6990 ± 0.0034	0.7832 ± 0.0027

Table 4: Text classification comparison results *w.r.t.* test accuracy (mean ± standard deviation). HyperGAT significantly outperforms all its variants on each dataset based on t-tests ( $p < 0.05$ ).

tion encoded by the sequential hyperedges is more important, but adding semantic hyperedges can enhance the model expressiveness. It also indicates that heterogeneous high-order context information can complement each other and we could investigate more meaningful hyperedges to further improve the performance of our approach.

#### 4.4 Case Study

**Embedding Visualization.** In order to show the superior embedding quality of HyperGAT over other methods, we use t-SNE (Maaten and Hinton, 2008) to visualize the learned representations of documents for comparison. Specifically, Figure 4 shows the visualization results of the best performing baseline Text-level GNN and HyperGAT on the test documents of Ohsumed. Note that the node’s color corresponds to its label, which is used to verify the model’s expressive power on 23 document classes. From the embedding visualization, we are able to observe that HyperGAT can learn more expressive document representations over the state-of-the-art method Text-level GNN.

**Attention Visualization.** To better illustrate the learning process of the proposed dual attention mechanism, we take a text document from 20NG (labeled as *sport.baseball* correctly) and visualize the attention weights computed for the word *player*. As shown in Figure 5, *player* is con-

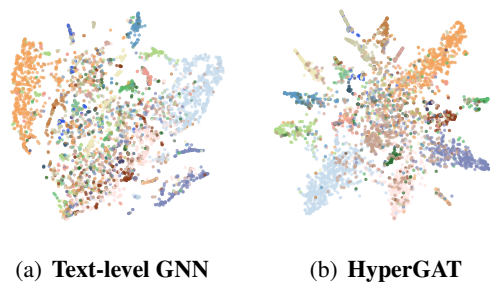


Figure 4: The t-SNE visualization of Text-level GNN and HyperGAT for test documents in Ohsumed.

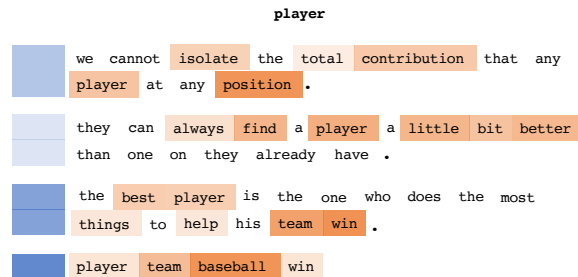


Figure 5: Visualization of the dual attention mechanism in HyperGAT. Figure best viewed in color.

nected to four hyperedges within the constructed document-level hypergraph. The first three lines ended with periods represent sequential hyperedges, while the last one without a period is a semantic hyperedge. Note that we use orange to denote the node-level attention weight and blue to denote the edge-level attention weight. Darker color represents larger attention weight.

On the one hand, node-level attention is able to select those nodes (words) carrying informative context on the same hyperedge. For example, *win* and *team* in the third hyperedge gain larger attention weights since they are more expressive compared to other words in the same sentence. On the other hand, edge-level attention can also assign fine-grained weights to highlight meaningful hyperedges. As we can see, the last hyperedge that connects *player* with *baseball* and *win* receives higher attention weight since it can better characterize the meaning of *player* in the document. To summarize, this case study shows that our proposed dual attention can capture key information at different granularities for learning expressive text representations.

## 5 Conclusion

In this study, we propose a new graph-based method for solving the problem of inductive text classification. Apart from the existing efforts, we propose to model text documents with document-



level hypergraphs and further develop a new family of GNN model named HyperGAT for learning discriminative text representations. Specifically, our method is able to acquire more expressive power with less computational consumption for text representation learning. By conducting extensive experiments, the results demonstrate the superiority of the proposed model over the state-of-the-art methods.

## Acknowledgements

This material is in part supported by the National Science Foundation (NSF) grant 1614576.

## References

- Song Bai, Feihu Zhang, and Philip HS Torr. 2020. Hypergraph convolution and hypergraph attention. *Pattern Recognition*.
- Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. 2018. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research*.
- Yiming Cui, Zhipeng Chen, Si Wei, Shijin Wang, Ting Liu, and Guoping Hu. 2017. Attention-over-attention neural networks for reading comprehension. *ACL*.
- Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *NeurIPS*.
- Kaize Ding, Jundong Li, Rohit Bhanushali, and Huan Liu. 2019a. Deep anomaly detection on attributed networks. In *SDM*.
- Kaize Ding, Yichuan Li, Jundong Li, Chenghao Liu, and Huan Liu. 2019b. Feature interaction-aware graph neural networks. *arXiv preprint arXiv:1908.07110*.
- Kaize Ding, Jianling Wang, Jundong Li, Kai Shu, Chenghao Liu, and Huan Liu. 2020. Graph prototypical networks for few-shot learning on attributed networks. *arXiv preprint arXiv:2006.12739*.
- Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. 2019. Hypergraph neural networks. In *AAAI*.
- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. 2017. Neural message passing for quantum chemistry. In *ICML*.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NeurIPS*.
- Lianzhe Huang, Dehong Ma, Sujian Li, Xiaodong Zhang, and Houfeng WANG. 2019. Text level graph neural network for text classification. *EMNLP*.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. *NeurIPS*.
- Hu Linmei, Tianchi Yang, Chuan Shi, Houye Ji, and Xiaoli Li. 2019. Heterogeneous graph attention networks for semi-supervised short text classification. In *EMNLP*.
- Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. 2016. Recurrent neural network for text classification with multi-task learning. *IJCAI*.
- Xien Liu, Xinxin You, Xiao Zhang, Ji Wu, and Ping Lv. 2020. Tensor graph convolutional networks for text classification. In *AAAI*.
- Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*.
- Riccardo Miotto, Li Li, Brian A Kidd, and Joel T Dudley. 2016. Deep patient: an unsupervised representation to predict the future of patients from the electronic health records. *Scientific reports*.
- Dinghan Shen, Guoyin Wang, Wenlin Wang, Martin Renqiang Min, Qinliang Su, Yizhe Zhang, Chunyuan Li, Ricardo Henao, and Lawrence Carin. 2018. Baseline needs more love: On simple word-embedding-based models and associated pooling mechanisms. In *ACL*.
- Kai Sheng Tai, Richard Socher, and Christopher D Manning. 2015. Improved semantic representations from tree-structured long short-term memory networks. *ACL*.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. *ICLR*.

- Guoyin Wang, Chunyuan Li, Wenlin Wang, Yizhe Zhang, Dinghan Shen, Xinyuan Zhang, Ricardo Henao, and Lawrence Carin. 2018. Joint embedding of words and labels for text classification. In *ACL*.
- Jianling Wang, Kaize Ding, Liangjie Hong, Huan Liu, and James Caverlee. 2020. Next-item recommendation with sequential hypergraphs. In *WSDM*.
- Sida Wang and Christopher D Manning. 2012. Baselines and bigrams: Simple, good sentiment and topic classification. In *ACL*.
- Yequan Wang, Minlie Huang, Xiaoyan Zhu, and Li Zhao. 2016. Attention-based lstm for aspect-level sentiment classification. In *EMNLP*.
- Felix Wu, Tianyi Zhang, Amauri Holanda de Souza Jr, Christopher Fifty, Tao Yu, and Kilian Q Weinberger. 2019a. Simplifying graph convolutional networks. *ICML*.
- Man Wu, Shirui Pan, Xingquan Zhu, Chuan Zhou, and Lei Pan. 2019b. Domain-adversarial graph neural networks for text classification. In *ICDM*.
- Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. Hierarchical attention networks for document classification. In *NAACL-HLT*.
- Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Graph convolutional networks for text classification. In *AAAI*.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *NeurIPS*.
- Jie Zhou, Ganqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2018. Graph neural networks: A review of methods and applications. *arXiv preprint arXiv:1812.08434*.

## A Appendix

### A.1 Implementation Details

As the supplement to Section 4, in the following, we explain the implementation of HyperGAT.

**LDA Model Training.** We use the implementation provided in scikit-learn to train the LDA model. We only use the documents in the training set to train the LDA model for each dataset. We select to use the Online Variational Bayes method for model learning. We set the random state is set to be 0 and the learning offset to be 50. As for the other parameters, we follow the default setting provided by scikit-learn. The topic number is set to be the same as the number of classes for each of

the datasets. And we select the Top-10 keywords of each topic to construct the semantic hyperedges.

**Implementation of HyperGAT.** The proposed HyperGAT model is implemented in PyTorch and optimized with the Adam optimizer (Kingma and Ba, 2014). It is trained and tested on a 12 GB Titan Xp GPU. Specifically, the hypergraph attention network consists of two layers with 300 and 100 embedding dimensions, respectively. We use one-hot vectors as the node attributes. The batch size is set to be 8 for all the datasets. We grid search for the learning rate in  $\{0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1\}$ , L2 regularization in  $\{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$  and the dropout rate in  $\{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7\}$ . The optimal values are selected when the model achieves the highest accuracy for the validation samples. The optimized learning rate  $\alpha$  for MR is 0.0005 while that for the other datasets is 0.001. We select the L2 regularization to be  $10^{-6}$  and the dropout rate to be 0.3 for the best performance. For each dataset, we train the model for 100 epochs or stop if the performance for the validation doesn't increase for 5 consecutive epochs as an early-stopping strategy. Under the optimized setup, the model can converge in 587s, 145s, 156s, 97s and 78s on average for 20NG, R8, R52, Ohsumed and MR, respectively.

**Validation Performance.** As supplement to the test results in Table 2, we also report the corresponding validation performance of the proposed HyperGAT. The validation accuracy is  $0.9355 \pm 0.0011$ ,  $0.9755 \pm 0.0019$ ,  $0.9375 \pm 0.0023$ ,  $0.6964 \pm 0.0024$  and  $0.7779 \pm 0.0015$  for 20NG, R8, R52, Ohsumed and MR, respectively.

### A.2 Space Complexity Analysis

Theoretically, the main difference of memory usage between HyperGAT and other methods lies in the size of the adjacency matrix. Formally, let  $N$  denote vocabulary size and  $M$  denote document size. Take TextGCN as an example, the size of the adjacency matrix is  $(N + M)^2$ . As HyperGAT adopts document-level hypergraphs, for each hypergraph, the adjacency matrix size is  $n \times m$ , where  $n$  is the number of words and  $m$  is the number of hyperedges in a document. Based on mini-batch training, the memory consumption for each mini-batch is about  $n \times m \times bsz$ . Since  $N$  and  $M$  are way larger than  $n$  and  $m$ , the memory consumption of HyperGAT can be largely reduced in practice.