# CONTRASTIVE LEARNING WITH SIMPLICIAL CONVO LUTIONAL NETWORKS FOR SHORT-TEXT CLASSIFICA TION

#### Anonymous authors

Paper under double-blind review

#### ABSTRACT

Text classification is a fundamental task in Natural Language Processing (NLP). Short text classification has recently captured much attention due to its increased amount from various sources with limited labels and its inherent challenges for its sparsity in words and semantics. Recent studies have adopted self-supervised contrastive learning across different representations to improve performance. However, most of the current models face several challenges. Firstly, the augmentation step might not be able to generate positive and negative samples that are semantically similar and dissimilar to the anchor respectively. Secondly, the text data could be enhanced with external auxiliary information that might introduce noise to the sparse text data. In addition, they are limited in capturing higher-order information such as group-wise interactions. In this work, we propose a novel document simplicial complex construction based on text data for a higher-order message-passing mechanism. We develop a simplicial complex representation for text sentences based on the directed word co-occurrence. Novel features are proposed for 0-simplex (word), 1-simplex (word-pair), and 2-simplex (three consecutive words) to characterise intrinsic higher-order structural information among words. We also enhance the short text classification performance by contrasting the structural representation with the sequential representation generated by the transformer mechanism for improved outcomes. The proposed framework, Contrastive Learning with Simplicial Convolutional Networks (C-SCN), leverages the expressive power of graph neural networks, models higher-order information beyond pair-wise relations and enriches features through contrastive learning. Experimental results on four benchmark datasets demonstrate the capability of C-SCN to outperform existing models in analysing sequential and complex shorttext data.

030

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

028

029

031

032

034

#### 1 INTRODUCTION

Text classification is a fundamental task in Natural Language Processing (NLP). It involves 040 analysing the content of texts and determining which predefined category they belong to based on 041 their representation. Unlike longer texts, short texts have recently captured much attention, with 042 an increase in the number appearing in various sources, such as social media, search snippets, and 043 news feeds. However, these short texts with a few words pose challenges to the current models in 044 generating effective representations and are not usually labelled in real-world cases (Linmei et al. 045 (2019)). Supervised learning on short-text classification has gained significant attention and has 046 been applied to different tasks for web reviews (Pang & Lee (2004)), news feed (Yao et al. (2019)) 047 and medical information (Liu et al. (2020)). On the other hand, data labelling has been expensive, 048 labour-intensive and time-consuming. Few-shot learning has been popular with low resources required by training on a few labelled samples, either with or without pre-training. Furthermore, graph models have been widely used to capture complex relationships between text data's structural, se-051 mantic, and syntactic meanings. To address the label scarcity issue, contrastive learning has been adopted to enhance performance. Many researchers (Sun et al. (2022); Wen & Fang (2023); Liu 052 et al. (2024)) have explored the effectiveness of combining graph models and contrastive learning within the scope of few-shot learning.

1

054 Although they have achieved successful outcomes, some limitations and challenges still exist. 055 Firstly, the data augmenting step and the negative sampling step of contrastive learning might distort 056 the semantic meaning and introduce unnecessary noise. For example, removing graph components 057 is adopted as a data augmentation strategy, but it might disrupt the text's original meaning. An in-058 stance from the Movie Review (MR) dataset (Pang & Lee (2005)): "there's not enough to sustain the comedy" while removing the word "not" reversely changes the meaning of this short sentence. Furthermore, negative sampling of texts with different syntaxes while similar semantics might be 060 designed to be pushed away from each other. Secondly, some auxiliary information such as enti-061 ties, latent topics, and part-of-speech (POS) tags (such as nouns and verbs) might be added to graph 062 models for language understanding and enriching the limited available local context. However, this 063 step might introduce misinformation, such as pulling documents that express opposite semantics but 064 similar topics closer. Lastly, graph models are mathematically limited in modelling higher-order fea-065 tures, such as group-wise interactions among a few nodes and edges expressed in terms of phrases. 066 For example, the short sentence "It is what it is" uses repetition to emphasise the acceptance of 067 the status quo. At the same time, graph models with only nodes and edges learn pairwise interac-068 tion. They need to extend the number of layers in order for words to incorporate the meaning of 069 other words further apart. Group-wise phrase "it is" needs to be linked with "what" to model such repetition.

071 To address the challenges mentioned above, a novel model combining higher-order features with 072 contrastive learning is proposed in this paper called Contrastive Learning with Simplicial Convo-073 lutional Networks (C-SCN) for short-text classification tasks. Specifically, SCN adopts simplicial 074 complex to robustly model richer and more complex information for better document understanding. 075 Document simplicial complexes are firstly constructed based on text data, and respective features are defined for simplexes of different complexities. We further integrate the features into an inductive 076 message-passing mechanism, considering long-range structural information for individual document 077 simplicial complexes. Furthermore, contrastive learning is embraced to compare the structural representation from SCN and sequential representation from the transformer model so that the power 079 of both sides can be combined for better performance in a few-shot learning setting. 080

Our work's key contributions are as follows. Firstly, we propose the construction of simplicial complexes based on text data and define features on 0-simplexes, 1-simplexes and 2-simplexes in the context of short-text classification in the message-passing mechanism of SCN. Secondly, we extend SCN with contrastive learning, such as C-SCN, where the power of sequential representation from the transformer model is integrated to solve the existing limitations and challenges. Lastly, the experiment with C-SCN in benchmark short text classification tasks demonstrates better results than competitive baseline models in the few-shot setting.

The remainder of the paper is organised as follows: Section 2 reviews the literature on graph neural networks, contrastive learning, and neural networks on simplexes. Section 3 outlines the proposed model structure for message passing on higher-order structures and contextualises the methods in text classification. Section 4 introduces four short text classification task datasets from various domains used in the experiments. Section 5 presents the performance metrics compared with other models and ablation studies. Finally, Section 6 concludes the work and proposes future directions.

- 2 LITERATURE REVIEW
- 095 096

094

### 2.1 GRAPH NEURAL NETWORKS IN SHORT-TEXT CLASSIFICATION

098 Graph neural networks (GNN) are powerful deep learning models to model representations of struc-099 tural data (Scarselli et al. (2009)). Through a message-passing mechanism, features of nodes and 100 edges are aggregated in the neighbourhood formed by components in the local document. Texts 101 could be used to construct different types of graphs, such as heterogeneous graphs (Yao et al. (2019)), 102 knowledge graphs (Ye et al. (2019)), dynamic graphs (Chen et al. (2020)), and hypergraphs (Ding 103 et al. (2020)). Early graph neural networks, such as Graph Convolutional Networks (GCN) (Kipf 104 & Welling (2017)), Graph Isomorphism Networks (GIN) (Xu et al. (2019)) and Graph Attention 105 Networks (GAT) (Veličković et al. (2018)), are integrated with the text graphs for improved results. Recently, model fusion has been adopted by Lin et al. (2021) to jointly train the transformer model 106 BERT with graph models with text data. On the other hand, graph models are limited in modelling 107 higher-order information in the group-wise form.

## 108 2.2 CONTRASTIVE LEARNING

110 Unlike traditional supervised learning, where a label is required for training, contrastive learning is a 111 self-supervised technique where augmented views of the same object are used to train a model which could gather positive samples closer and negative samples further apart (Jaiswal et al. (2021)). With 112 downstream tasks, contrastive learning has been actively applied in short-text classification scenarios 113 with few-shot settings. Sun et al. (2022) integrates the heterogeneous graph attention mechanism 114 with neighbouring contrastive learning to enrich the terms beyond the document and extend the 115 relations among documents; Wen & Fang (2023) pre-trains text and graph encoders followed by 116 few-shot and zero-shot fine-tuning process; Liu et al. (2024) innovates in augmented view of graph 117 features through the singular value decomposition (SVD) of the feature matrix and in assigning 118 weak labels to document through k-means clustering. On the other hand, these current methods 119 require a large amount of resources in preprocessing in the form of pre-training or enriching text 120 with additional information, such as entity recognition and POS tagging processes. The augmented 121 view in contrastive learning might also introduce unnecessary noise and misleading information.

122

# 123 2.3 TOPOLOGICAL DEEP LEARNING (TDL)

125 Topological deep learning combines the techniques from deep learning and topological tools that structure data manifolds (Zia et al. (2024)). Topological representations, including cell complexes 126 (Hajij et al. (2020); Giusti et al. (2023); Bodnar et al. (2022)), simplicial complexes (Bodnar (2022); 127 Schaub et al. (2022)), combinatorial complexes (Hajij et al. (2023)), sheaves (Hansen & Ghrist 128 (2019)) and hypergraphs (Feng et al. (2018); Bai et al. (2021)), model not only pair-wise interac-129 tions that are present on a graph, but also higher-order interactions among three elements or more. 130 Algebraic topology-based methods have achieved noteworthy results in protein analysis (Xia & 131 Wei (2014); Sverrisson et al. (2021); Wee & Xia (2022)), virus analysis (Chen et al. (2022)), drug 132 design (Cang & Wei (2017)) and material property classification (Reiser et al. (2022); Townsend 133 et al. (2020)), where topological representations demonstrate their robustness against deformation 134 and noise. Extending from algebraic topology-based methods, TDL employs the message-passing 135 mechanism on higher-order components, where the communication of information has been propa-136 gated through any neighbourhood relations (Roddenberry et al. (2021); Bodnar (2022); Hajij et al. 137 (2023)). However, there is a lack of studies on non-time-series sequential analysis with TDL on the text data, and we aim to explore its usage in the new field. 138

139 140

141 142

143

#### 3 Methods

3.1 SIMPLICIAL CONVOLUTIONAL NETWORKS (SCN)

We first provide the necessary details related to constructing document simplicial complexes, fol lowed by the message-passing mechanism on the higher-order structures.

Denote  $\mathcal{K}_k$  the set of k-simplexes for  $\mathcal{K}$ .  $\mathcal{K}_0$  will be referred to as the set of **0-simplexes** (nodes). 152  $\mathcal{K}_1$  refers to the set of **1-simplexes** (edges) and  $\mathcal{K}_2$  refers to the set of **2-simplexes** ("filled" trian-153 gles). For the text classification task, we construct the document as a simplicial complex with initial 154 representations of 0-simplexes, 1-simplexes, and 2-simplexes, as shown in Figure 1. We embrace 155 the bag-of-word model (Harris (1954)) and treat each word and punctuation as distinct 0-simplexes 156 initialised from GloVe embeddings (Pennington et al. (2014)). The three types of direction of 1-157 simplexes follow the sequential order of the tokens in each text as shown in Figure 1. 2-simplexes 158 are formed when any three words form a "filled" triangle. We differentiate their nine identities by 159 the neighbouring 1-simplexes for the 2-simplexes to be formed. An example of the 1-simplex  $e_{st}$  is shown in Figure 2 where the types of 2-simplex formed are determined by the 1-simplex between 160 s and o and the 1-simplex between t and o. It is to be noted that the self-loop is not considered 161 part of the 2-simplex formation since we consider unique 0-simplexes appearing in texts. One target



Figure 1: One document simplicial complex constructed for a document example from the Snippets dataset (Phan et al. (2008)) with different types of flow directions. Words and punctuation are tokenised into individual *0-simplexes* (nodes). *1-simplexes* (edges) are formed if 0-simplexes are next to each other with directions in chronological order. Lastly, *2-simplexes* (triangles) are constructed if the three words form a "filled" triangle. Three types of 1-simplexes are illustrated: (1) **Forward 1-simplexes** are the ones following the chronological order which points to the word that first appears in the text; (2) **Backward 1-simplexes** are 1-simplexes pointing to the word which is used before and referenced again; (3) **Self-loop 1-simplexes** are formed when 1-simplexes connect the same word.



Figure 2: 2-simplexes types for a 1-simplex with source 0-simplex s and target 0-simplex t. For the 1-simplex  $e_{st}$  with a source 0-simplex s, a target 0-simplex t and the defined direction from s to t, nine types of 2-simplexes could define the information flow through the 2-simplex. For a 0simplex o that forms a 2-simplex with the target 1-simplex  $e_{st}$ , there exist three types of 1-simplexes between 0-simplex o and 0-simplex s, as well as between 0-simplex o and 0-simplex t: into, out of or bidirectional, resulting in nine types of 2-simplex with respect to the 1-simplex  $e_{st}$ .

193 194

162

163

164

167

169

170

171 172

173 174

183

185 186 187

188

189

190

191

192

1-simplex could be part of multiple 2-simplexes. The 1-simplex and 2-simplex embeddings will be
 initialised randomly, and all embedding matrices are optimised during training.

197 The message-passing mechanism leverages the connectivity information in simplicial complexes. 198 For a simplex  $\sigma_k$ , we denote its **boundary adjacent simplexes**  $\mathcal{B}(\sigma_k)$  as the set of lower-199 dimensional simplexes  $\sigma_{k-1}$  on the boundary of  $\sigma_k$ , its **co-boundary adjacent simplexes**  $\mathcal{C}(\sigma_k)$ 200 as the set of higher-dimensional simplexes  $\sigma_{k+1}$  with  $\sigma_k$  on their boundaries, its **lower adjacent** 201 **simplexes**  $\mathcal{N}_{\downarrow}(\sigma_k)$  as those with the same dimension as  $\sigma_k$  that share a lower-dimensional simplex 202  $\sigma_{k-1}$  on their boundary, and its **upper adjacent simplexes**  $\mathcal{N}_{\uparrow}(\sigma_k)$  as those with the same dimen-203 sion as  $\sigma_k$  that are on the boundary of the same higher-dimensional simplex  $\sigma_{k+1}$  (Bodnar et al. 204 (2021)).

205 The hidden representation of simplexes will be initialised with vectors  $\mathbf{x}_{\sigma_i}$  for all  $\sigma_k \in \mathcal{K}_k, k \in$ 206  $\{0, 1, \dots, K\}$  and we set K = 2. That means at the initial state, the layer representation for simplex  $\sigma_k$  is  $h_{\sigma_k}^{(0)} = \mathbf{x}_{\sigma_k}$  for  $k \in \{0, 1, \dots, K\}$ . The messages are then aggregated according to the neighbourhood in which the simplexes sit: at the state  $\ell + 1$  and for the target k-simplex  $\sigma_k$ , the 207 208 209 message function  $M_k^{(\ell+1)}$  collects information from neighbouring simplexes of the same dimension 210  $\sigma'_k \in \mathcal{N}(\sigma_k)$  where  $\mathcal{N}(\sigma_k) = \mathcal{N}_{\downarrow}(\sigma_k) \cup \mathcal{N}_{\uparrow}(\sigma_k)$ , those of one dimension lower  $\sigma_{k-1} \in \mathcal{B}(\sigma_k)$ 211 and those of one dimension higher  $\sigma_{k+1} \in \mathcal{C}(\sigma_k)$  as illustrated in Equation (1). In the context 212 of text classification, we adopt the message function as a multi-layer perception (MLP) for both 213 0-simplexes  $\sigma_0 \in \mathcal{K}_0$  and 1-simplex  $\sigma_1 \in \mathcal{K}_1$  updates. The message passing is set to collect information from neighbouring simplexes and co-boundary adjacent simplexes. For the aggregation 214 of all the components in the document simplicial complex, we adopt row summation as illustrated 215 in Equation (2).

=

For individual  $\sigma_k \in \mathcal{K}_k$  and  $k \in \{0, 1, \cdots, K\}$ ,

218 219

$$m_{\sigma_{k}}^{(\ell+1)} = \underset{\substack{\sigma_{k}' \in \mathcal{N}(\sigma_{k}) \\ \sigma_{k-1} \in \mathcal{B}(\sigma_{k}) \\ \sigma_{k+1} \in \mathcal{C}(\sigma_{k})}}{\operatorname{AGG}} \left( \phi \left( M_{k}^{(\ell+1)}(h_{\sigma_{k}'}^{(\ell)}, h_{\sigma_{k-1}}^{(\ell)}, h_{\sigma_{k+1}}^{(\ell)}) \right) \right)$$
(1)

221 222 223

220

$$= \sum_{\substack{\sigma'_k \in \mathcal{N}(\sigma_k) \\ \sigma_{k+1} \in \mathcal{C}(\sigma_k) \cap \mathcal{C}(\sigma'_k)}} \phi\left(\mathsf{MLP}_k^{(\ell+1)}\left(h_{\sigma'_k}^{(\ell)} + h_{\sigma_{k+1}}^{(\ell)}\right)\right)$$
(2)

225 226 227

228

229

224

where for k = 0, we do not consider simplex of dimension (n - 1),  $h_{\sigma_k}^{(\ell)}$  refers to the simplex  $\sigma_k$ 's feature at the state  $\ell$ .  $h_{\sigma_{k+1}}^{(\ell)}$  is set to a zero vector if  $\mathcal{C}(\sigma_k) \cap \mathcal{C}(\sigma'_k)$  is empty. MLP<sup> $(\ell+1)$ </sup> refers to trainable multi-layer perception at the state  $\ell + 1$ .

Similarly to the GNN framework, the update function UPDATE<sup> $(\ell+1)$ </sup> will synchronise the representation of the *k*-simplex to the new state, as shown in Equation (3), and we adopt the Gated Recurrent Unit (GRU) for the text classification task.

 $h_{\sigma_{k}}^{(\ell+1)} = \text{UPDATE}_{k}^{(\ell+1)}(h_{\sigma_{k}}^{(\ell)}, m_{\sigma_{k}}^{(\ell+1)}) = \text{GRU}_{k}^{(\ell+1)}(h_{\sigma_{k}}^{(\ell)}, m_{\sigma_{k}}^{(\ell+1)})$ (3)

236 Lastly, the readout function READOUT will obtain the representation for the document simplicial 237 complex by pooling k-simplexes' features of the final state L in Equation (4). A global self-attention 238 mechanism (Lin et al. (2017)) is specifically applied for text data, summarising the 0-simplexes 239 and 1-simplexes. For the final layer representation  $h_{\mathcal{K}}^L$  of the document simplicial complex with 240 0-simplexes  $\sigma_0 \in \mathcal{K}_0$  and 1-simplexes  $\sigma_1 \in \mathcal{K}_1$ , its individual simplex attention score  $\alpha_{\sigma_k}$  is derived with two multi-layer perceptions without bias denoted by  $W_1$  and  $W_2$ . The final simplex 241 representation for the document simplicial complex,  $h_{\mathcal{K}}^L$ , is hence the summation of the attention 242 score multiplied by the respective final simplex features  $h_{\sigma_k}^L$  for  $k \in \{0, 1\}$ . 243

$$h_{\mathcal{K}}^{L} = \operatorname{READOUT}\left(\{h_{k}^{(L)} | k \in \{0, 1, \cdots, K\}\}\right)$$
(4)

$$= \left(\sum_{\sigma_0 \in \mathcal{K}_0} \alpha_{\sigma_0} h_{\sigma_0}^L\right) \oplus \left(\sum_{\sigma_1 \in \mathcal{P}_1} \alpha_{\sigma_1} h_{\sigma_1}^L\right)$$
(5)

$$\alpha_{\sigma_k} = \frac{\exp\left(\tanh(\mathbf{W}_1 h_{\sigma_k}^{-}) \cdot \mathbf{W}_2\right)}{\sum_{\sigma'_k \in \mathcal{K}_0} \exp\left(\tanh(\mathbf{W}_1 h_{\sigma'_k}^L) \cdot \mathbf{W}_2\right)}$$
(6)

252 253

251

254 where  $h_k^{(L)}$  refers to the final collective representation for all  $\sigma_k \in \mathcal{K}_k$ . Finally, a linear layer 255 with a softmax classifier will transform the results to the same dimension as the label set and make 256 predictions. A summary of the proposed SCN framework is illustrated in Figure 3. With the input 257 sentence "a thriller without a lot of thrills.", a simplicial complex could be constructed with the 258 following components. For 0-simplexes, we have matches  $v_1$ : "a",  $v_2$ : "thriller",  $v_3$ : "without",  $v_4$ : 259 "lot",  $v_5$ : "of",  $v_6$ : "thrills", and  $v_7$ : ".". For 1-simplexes,  $e_1, e_2, e_4, e_5, e_6, e_7$  are forward edges 260 connecting two words, and  $e_3$  is a backward edge between "without" and "a". Lastly, a 2-simplex au is a type-4 triangle formed by the 1-simplexes connecting among the words "a", "thriller" and 261 "without". 262

Assuming message functions are fully connected neural networks, the SCN could be evaluated with three components: feature transformation in neural networks, neighbourhood aggregation and nonlinear activation. Assuming that all the layers are of the same size *F* and the embedding size is fixed with *F* for 0-simplexes, 1-simplexes and 2-simplexes, the features are initialised from all three kinds of simplexes and the dense matrix multiplication takes  $\mathcal{O}(|\mathcal{K}_0|F^2+|\mathcal{K}_1|F^2+|\mathcal{K}_2|F^2) = \mathcal{O}(|\mathcal{K}|F^2)$ . The aggregation and update step will take  $\mathcal{O}(|\mathcal{K}_1|F^2+|\mathcal{K}_2|F^2)$  for 0-simplex and 1-simplex updates. Non-linear activation is an element-wise function which will take  $\mathcal{O}(|\mathcal{K}_0|+|\mathcal{K}_1|)$ . As a result, over *L* layers, the final time complexity is  $\mathcal{O}(|\mathcal{K}_0|+|\mathcal{K}_1|+|\mathcal{K}_1|F^2+|\mathcal{K}_2|F^2) = \mathcal{O}(|\mathcal{K}|F^2)$ .



Figure 3: Message-passing mechanism in Simplicial Convolutional Networks (SCN) up to two-289 dimension. The above figure illustrates an example of a simplicial complex for a seven-token sen-290 tence with a message-passing mechanism that collects neighbouring information from the same di-291 mension, one dimension lower and one dimension higher. The input simplicial complex  $\mathcal{K}$  consists 292 of 0-simplexes  $v_1, \dots, v_7$ , 1-simplexes  $e_1, \dots, e_7$  and the 2-simplex  $\tau$ . Pre-defined and trainable 293 features of different simplexes are used as input on the left-hand side. SCN leverages neighbouring, boundary and co-boundary simplexes to carry feature information and update the target 0-simplexes 295 and 1-simplexes separately in different layer states. Finally, the features of different simplexes are 296 read out for downstream tasks. 297

#### 299 3.2 SCN with Contrastive Learning

To alleviate the abovementioned challenges with contrastive learning, we adopt a dual-encoder
 framework inspired by Wen & Fang (2023) where we generate text representations from transformer
 blocks and graph representations from SCN in parallel; hence, the training process could optimise
 the contrastive learning and classification task as shown in Figure 4.

We employ the BERT model (Devlin et al. (2019)) as the text encoder and SCN encoder that digest the document data  $\mathbf{x}_{doc}$ . We denote  $Z_t, Z_s \in \mathbb{R}^{\gamma}$  as the text encoder and SCN encoder output. MLP( $\bullet$ ) is a linear layer that processes the output to the target space's dimension  $\gamma$ .

$$Z_t = \mathrm{MLP}_t(\mathrm{BERT}(\mathbf{x}_{\mathrm{doc}})), \ Z_s = \mathrm{MLP}_{sc}(\mathrm{SCN}(\mathbf{x}_{\mathrm{doc}}))$$
(7)

The constrastive loss is derived by the cross-entropy loss (CE) between the normalised (norm) text encoder output and the normalised SCN encoder.

$$\mathcal{L}_{cl} = \operatorname{CE}(\operatorname{norm}(Z_t), \operatorname{norm}(Z_s)) \tag{8}$$

At the same time, we include the training objective against the ground-truth label y, which is a linear interpolation of the text encoder and SCN encoder after transformation to the same dimension as the label space ( $\tilde{\bullet}$ ) inspired by Lin et al. (2021).

$$Z = \frac{1}{2} \left( \text{softmax}(\tilde{Z}_s) + \text{softmax}(\tilde{Z}_t) \right)$$
(9)

317

298

308 309

313

$$\mathcal{L}_{label} = \operatorname{CE}(Z, y) \tag{10}$$

The final loss function is the integration of the contrastive loss with the classification loss.

$$\mathcal{L} = \mathcal{L}_{label} + \eta \cdot \mathcal{L}_{cl} \tag{11}$$



Figure 4: The Contrastive Learning with SCN (C-SCN) framework. The transformer encoder and SCN encoder will generate text representation  $Z_t$  and document simplicial complex representation  $Z_s$  respectively. The learned features will be used for contrastive learning by comparing themselves. Meanwhile, the two representations will contribute to the classification task with equal weights. As a result, the final loss  $\mathcal{L}$  includes the contrastive loss  $\mathcal{L}_{cl}$  and  $\mathcal{L}_{label}$ .

- 340 where  $\eta$  is a control parameter.

4 EXPERIMENTS

#### 4.1 DATASETS

The experiments are conducted on four datasets for short text classification tasks. The datasets are briefly introduced below, and a summary table is reported in Table 1. We adopt the same data preprocessing techniques as Wang et al. (2017) with slight modifications to include punctuation, keep the hashtag messages and add self-connection. Twitter (Bird et al. (2009)) is a binary clas-sification dataset for sentiments "positive" and "negative" collected by Natural Language Toolkit. MR (Pang & Lee (2005)) contains movie review documents from Rotten Tomato with binary sen-timent categories. Snippets (Phan et al. (2008)) contains Google web search text data with eight categories: "business", "computer", "health", "sports", "culture and art", "education and science", "engineering" and "politics and society". StackOverflow (Hamner et al. (2012)) contains question text from StackOverflow, and we choose the samples as Xu et al. (2017) for 20,000 questions with 20 categories.

Table 1: Summary	statistics	for text	Datasets.
------------------	------------	----------	-----------

	2			
Dataset	Twitter	MR	Snippets	StackOverflow
# Doc	10,000	10,662	12,340	20,000
# Train	40	40	160	400
Train ratio	0.40%	0.38%	1.30%	2%
# Tokens	12,229	18,337	29,422	11,161
Avg. Length	9.3	20.4	18.0	9.3
# Class	2	2	8	20
Avg. # 1-simplexes	21.87	37.79	30.23	17.24
Avg. # 2-simplexes	0.24	0.74	3.24	0.21

#### 4.2 BASELINE MODELS

We compare C-SCN with other various types of benchmark language models for short-text classification as reported by Liu et al. (2024).

Traditional Language Models: TF-IDF (Rajaraman & Ullman (2011)) refers to the term frequency inverse document frequency, and it measures the importance of word tokens to the document. The
 features generated are passed in a support vector machine (SVM) (Crammer & Singer (2002)) for the
 classification task. LDA (Blei et al. (2003)) refers to Latent Dirichlet Allocation and extracts latent
 topics from the text through probabilistic models. The features are trained with SVM for short text classification. PTE (Tang et al. (2015)) refers to Predictive Text Embedding, which utilises
 heterogeneous text networks for embeddings.

Machine Learning Models: CNN (Kim (2014)) refers to Convolutional Neural Networks with pre trained GloVe word embeddings (Pennington et al. (2014)). LSTM (Liu et al. (2016)), which
 refers to Long-Short Term Memory, is trained GloVe embeddings. BERT (Devlin et al. (2019)),
 which refers to the Bidirectional encoder representations from transformers and its modified version
 RoBERTa (Zhuang et al. (2021)) leverages pre-training through self-supervised learning and could
 be fine-tuned to specific downstream tasks.

384 Graph-based Language Models: TL-GNN (Huang et al. (2019)) refers to text-level GNN, which 385 adopts small windows for texts to focus on local features. TextGCN (Yao et al. (2019)), Text Graph 386 Convolutional Network, constructs individual text graphs with document nodes based on word co-387 occurrences and word-document relations. TextING (Zhang et al. (2020)) adopts individual text 388 graphs and inductively trains the model. HyperGAT (Ding et al. (2020)), Hypergraph Attention Networks enhances the expressive power of graphs on text classification by including high-order 389 information and reducing computational resources needed for training. STGCN (Ye et al. (2020)), 390 which refers to the short text graph convolutional network, integrates BERT and the bidirectional 391 LSTM in graph models to enhance performance on short texts. **DADGNN** Liu et al. (2021), Deep 392 Attention Diffusion Graph Neural Networks, applies attention diffusion and decoupling techniques 393 targeting some limitations of GNN such as oversmoothing and restricted receptive field. 394

395 Graph-based Models with external knowledge beyond documents or Contrastive Learning: STCKA (Chen et al. (2019)) refers to Short Text Classification with Knowledge-powered Attention, which 396 utilises attention mechanisms and entity conceptualisation to enhance text features. HGAT (Linmei 397 et al. (2019)) is known as Heterogeneous Graph Attention Networks, and its enhanced version in-398 corporates topic and entity beyond the texts for enriched graphs. SHINE (Wang et al. (2021)) is a 399 hierarchical heterogeneous graph representation learning method for short text classification which 400 executes entity and POS tagging for various types of node features. NC-HGAT (Sun et al. (2022)) 401 integrates HGAT with neighbouring contrastive learning. GIFT (Liu et al. (2024)) is the graph con-402 trastive learning for short text classification that employs SVD and k-means clustering methods in 403 contrastive learning.

404 405

406

#### 4.3 IMPLEMENTATION DETAILS

407 Following with few-shot setting for short text classification framework (Sun et al. (2022); Wen & Fang (2023); Liu et al. (2024)), from each category, 20 samples are selected randomly to form the 408 train set, another 20 samples are selected randomly to form the validation set, and the rest are in-409 cluded in the unseen test set. The 0-simplex embeddings are initialised with GloVe embeddings 410 Pennington et al. (2014). The embedding matrices for 1-simplexes and 2-simplexes are randomly 411 initialised and optimised to size 128. The learning rate is  $1 \times 10^{-4}$ , and the batch size is 128. A 412 dropout rate of 50% is implemented to reduce the complexity of the model and prevent overfitting 413 problems. The model is trained with the PyTorch Geometric<sup>1</sup> package for 100 epochs with early 414 stopping where the validation loss does not improve for ten epochs. The best weights are obtained 415 from the model with the best validation accuracy. Cross-entropy loss is used with an Adam opti-416 miser. The experiments are conducted ten times with NVIDIA RTX A6000 with 48GB of memory. 417 We compare the results with strong baseline models with ten iterations of different training, valida-418 tion and test sets. The average test accuracies and F1 scores are used for comparison.

419 420

421 422

423

#### 5 RESULTS AND DISCUSSION

#### 5.1 Results

The experiment results are reported in Table 2. Compared with other competitive models, C-SCN has achieved the best test accuracies and F1 scores, indicating the model's ability to capture sentiments and sequential information in text documents.

We attribute the better performance to the following analysis. Firstly, we adopt SCN, a higher-order framework extending the expressive power of GNN. Features assigned to 0-simplexes, 1-simplexes and 2-simplexes could better represent the sentence structure and are generalised well across different contexts. The involvement of 1-simplexes and 2-simplexes in the message-passing mechanism

<sup>431</sup> 

<sup>&</sup>lt;sup>1</sup>https://pytorch-geometric.readthedocs.io/en/latest/index.html

/35	are underlined.			in une pui			cora, an		
435	Dataset	Tw	itter	MR		Snippets		StackOverflow	
/137	Metrics	F1	Acc	F1	Acc	F1	Acc	F1	Acc
120	TF-IDF	53.62	52.46	54.29	48.13	64.70	59.17	59.19	59.06
430	LDA	54.34	53.97	54.40	48.39	62.54	56.4	60.19	59.52
439	PTE	54.24	53.17	55.02	52.62	63.10	59.11	62.56	61.32
440		57.29	56.02	59.06	59.01	77.09	69.28	63.75	61.21
441	LSTM	60.28	60.22	60.89	60.70	75.89	67.72	61.62	60.49
442	BERT	54.92	51.16	51.69	50.65	79.31	78.47	66.94	67.26
443	RoBERTa	56.02	52.29	52.55	51.30	79.55	79.02	69.91	70.35
444	TL-GNN	59.02	54.56	59.22	59.36	70.25	63.29	62.09	61.91
445	TextGCN	60.15	59.82	59.12	58.98	77.82	71.95	67.02	66.51
446	TextING	59.62	59.22	58.89	58.76	71.10	70.65	65.37	64.63
447	HyperGAT	59.15	55.19	58.65	58.62	70.89	63.42	63.25	62.10
448	DADGNN	59.51	55.32	58.92	58.86	71.65	70.66	66.26	65.10
449	STCKA	57.56	57.02	53.25	51.19	68.96	61.27	59.72	59.65
450	STGCN	64.33	64.29	58.25	58.22	70.01	69.93	69.23	69.10
451	-	† <u>6</u> 3.21 <sup>-</sup>	57.02	62.75	62.36	82.36	74.44	67.35	66.92
452	SHINE	72.54	72.19	64.58	63.89	82.39	81.62	73.05	72.73
453	NC-HGAT	63.76	62.94	62.46	62.14	82.42	74.62	67.59	67.02
454	GIFT	73.16	73.16	65.21	<u>65.16</u>	83.73	82.35	83.07	<u>82.94</u>
455	SCN	66.13	67.25	61.15	61.93	76.13	75.66	76.85	74.04
456	C-SCN	75.61	76.09	69.46	69.87	84.97	85.56	84.15	83.87
457		1				1		1	

Table 2: Results of test accuracy (%) and test F1-score (%) for short text classification where the best results based on 95% confidence in the pairwise *t*-tests are in bold, and the second-best results are underlined.

432

459 also expands the receptive fields of individual 0-simplexes where long-range information can be 460 transmitted through shallow neural network layers, thereby enhancing the impact of 0-simplexes 461 on the entire document. The self-attentive readout function connects 0-simplexes and 1-simplexes, creating expressive document-level summaries. This has promoted the SCN to perform the best in 462 the benchmark datasets among the graph-based models without external information or contrastive 463 learning. Secondly, the contrastive learning framework allows C-SCN to capture both structural 464 and textual information in the few-shot setting. Both structural representation and sequential repre-465 sentation are treated as augmented views of each other. This has contributed to preventing helpful 466 information from being removed, avoiding introducing noise or external information and combining 467 the capabilities of both models. 468

In addition, we see that the large language models, such as BERT and RoBERTa, which leverage 469 numerous pre-training, are not performing favourably with a few available labels. In contrast, graph-470 based models with external auxiliary knowledge or contrastive learning, including HGAT, SHINE, 471 NC-HGAT, and GIFT, could achieve competitive results. External auxiliary knowledge, such as 472 entity recognition and POS tagging, might help enrich the semantic and syntactic meaning of the 473 original text. Still, it might be introducing extra noise and unnecessary messages to the text data, 474 as shown in the deterioration of results from STGCN to HGAT. Furthermore, contrastive learning 475 with perturbation of the graphs might inject misinformation about the text's meaning, explaining the 476 difference between NC-HGAT and GIFT. Furthermore, introducing the global network within the 477 small train set where the connectivity or clustering effect is explored might not be significant. This 478 could explain why our model could outperform SHINE and GIFT.

479

481

#### 480 5.2 ABLATION STUDIES

Ablation studies are conducted to remove individual components to verify the capability of higher order simplexes and contrastive learning in enhancing text understanding. The results are reported
 in Table 3. It is observed that removing contrastive learning deteriorates the results for both SCN
 and BERT. Regarding higher-order simplexes, the removal of any component might deprecate the
 test accuracies and F1 scores across all datasets. Moreover, the inclusion of 1-simplexes followed



Figure 5: Hyperparameter  $\eta$  sensitivity across different datasets.

by the inclusion of 2-simplexes improves the results respectively, highlighting the importance of higher-order simplicial complexes in document understanding.

Table 3: Results of test accuracy for ablation studies. "C-SCN - 0-simplex" means the 1-simplexes and 2-simplexes are both removed in the model, whereas "C-SCN - 1-simplex" refers to the removal of 2-simplexes from the model.

503	of 2-simplexes from the	model.							
504	Dataset	Twi	itter	M	R	Snip	pets	StackO	verflow
505	Metrics	F1	Acc	F1	Acc	F1	Acc	F1	Acc
506	BERT	54.92	51.16	51.69	50.65	79.31	78.47	66.94	67.26
507	SCN	66.13	67.25	61.15	61.93	76.13	75.66	76.85	74.04
508	C-SCN - 0-simplex	74.50	74.78	67.48	68.27	84.58	85.13	82.79	82.36
509	C-SCN - 1-simplex	74.91	75.41	68.54	68.77	84.75	85.32	83.08	82.58
510	C-SCN	75.61	76.09	69.46	<b>69.87</b>	84.97	85.56	84.15	83.87

The hyperparameter sensitivity of  $\eta$  is investigated across different datasets, and the results are visualised in Figure 5. The control parameter  $\eta$  indicates the weights of contrastive loss in the model 513 training process. We could observe that there are various types of impact on test performance. In general, the test performance varies between the value 0 (no contrastive loss) and 1 (higher weight 515 of contrastive loss), while  $\eta = 1$  results in lower performance compared to the case of no contrastive 516 loss. One explanation for such variation could be the need to balance the focus between achieving the correct label and synchronising model weights between SCN and the transformer model. In our 518 experiments, a grid search is conducted for the best performance for the best  $\eta$  values.

519 520 521

517

495 496 497

498

499 500

501

502

511

512

514

#### 6 CONCLUSION

522 In conclusion, we propose Contrastive Learning with Simplicial Convolutional Networks (C-SCN), 523 which incorporates higher-order information for sequence analysis and is applied in short text clas-524 sification tasks. The model constructs document simplicial complexes and develops a convolutional 525 network to incorporate the higher-order simplexes' message passing with a self-attention readout. 526 Furthermore, we integrate the transformer model to generate augmented views in the contrastive 527 learning framework. Extensive experiments that simulate the lack of label situation in a few-shot set-528 ting indicate that our model leverages advantages from both structural and sequential representation, learns long-range information and enhances textual understanding with contextualised 1-simplexes 529 and 2-simplexes during training. 530

531 In the future, we would like to explore the interpretability of higher-order simplexes and their roles 532 in text understanding. The impact of the number of 1-simplexes and 2-simplexes on the performance 533 of C-SCN is also worth attention, and it could be more inspected within the context of longer doc-534 uments. Leveraging SCN's expressiveness in sequential analysis could have more applications in 535 other fields, such as recommender systems and process mining.

536

#### REFERENCES 538

Song Bai, Feihu Zhang, and Philip H.S. Torr. Hypergraph convolution and hypergraph attention. Pattern Recognition, 110:107637, 2021. ISSN 0031-3203.

540 Steven Bird, Edward Loper, and Ewan Klein. Natural Language Processing with Python. O'Reilly 541 Media Inc., 2009. 542 David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. J. Mach. Learn. 543 Res., 3(null):993-1022, mar 2003. 544 Cristian Bodnar. Topological Deep Learning: Graphs, Complexes, Sheaves. PhD thesis, Apollo -546 University of Cambridge Repository, 2022. 547 548 Cristian Bodnar, Fabrizio Frasca, Nina Otter, Yu Guang Wang, Pietro Liò, Guido Montúfar, and Michael Bronstein. Weisfeiler and lehman go cellular: Cw networks. In Advances in Neural 549 Information Processing Systems 34 - 35th Conference on Neural Information Processing Systems, 550 *NeurIPS 2021*, pp. 2625–2640. Neural information processing systems foundation, 2021. 551 552 Cristian Bodnar, Francesco Di Giovanni, Benjamin Paul Chamberlain, Pietro Lio, and Michael M. 553 Bronstein. Neural sheaf diffusion: A topological perspective on heterophily and oversmoothing 554 in GNNs. In Advances in Neural Information Processing Systems, 2022. 555 Zixuan Cang and Guo-Wei Wei. Topologynet: Topology based deep convolutional and multi-task 556 neural networks for biomolecular property predictions. PLoS Computational Biology, 13(7): e1005690, 2017. 558 559 Jiahui Chen, Yuchi Qiu, Rui Wang, and Guo-Wei Wei. Persistent laplacian projected omicron ba. 4 560 and ba. 5 to become new dominating variants. Computers in Biology and Medicine, 151:106262, 2022. 561 562 Jindong Chen, Yizhou Hu, Jingping Liu, Yanghua Xiao, and Haiyun Jiang. Deep short text clas-563 sification with knowledge powered attention. In Proceedings of the Thirty-Third AAAI Con-564 ference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelli-565 gence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, 566 AAAI'19/IAAI'19/EAAI'19. AAAI Press, 2019. 567 Yu Chen, Lingfei Wu, and Mohammed Zaki. Iterative deep graph learning for graph neural net-568 works: Better and robust node embeddings. In Advances in Neural Information Processing Sys-569 tems, volume 33, pp. 19314–19326. Curran Associates, Inc., 2020. 570 571 Koby Crammer and Yoram Singer. On the algorithmic implementation of multiclass kernel-based 572 vector machines. J. Mach. Learn. Res., 2:265-292, mar 2002. 573 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep 574 bidirectional transformers for language understanding. In Proceedings of the 2019 Conference 575 of the North American Chapter of the Association for Computational Linguistics: Human Lan-576 guage Technologies, NAACL-HLT 2019, pp. 4171–4186. Association for Computational Linguis-577 tics, 2019. 578 579 Kaize Ding, Jianling Wang, Jundong Li, Dingcheng Li, and Huan Liu. Be more with less: Hypergraph attention networks for inductive text classification. In Proceedings of the 2020 Conference 580 on Empirical Methods in Natural Language Processing (EMNLP), pp. 4927–4936, 2020. 581 582 Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. Hypergraph neural networks. 583 AAAI 2019, 2018. 584 585 Lorenzo Giusti, Claudio Battiloro, Lucia Testa, Paolo Di Lorenzo, Stefania Sardellitti, and Sergio Barbarossa. Cell attention networks. In 2023 International Joint Conference on Neural Networks 586 (IJCNN), pp. 1–8, 2023. 587 588 Mustafa Hajij, Kyle Istvan, and Ghada Zamzmi. Cell complex neural networks. In TDA & Beyond, 589 2020. 590 591 Mustafa Hajij, Ghada Zamzmi, Theodore Papamarkou, Nina Miolane, Aldo Guzmán-Sáenz, Karthikeyan Natesan Ramamurthy, Tolga Birdal, Tamal K. Dey, Soham Mukherjee, Shreyas N. 592 Samaga, Neal Livesay, Robin Walters, Paul Rosen, and Michael T. Schaub. Topological deep learning: Going beyond graph data, 2023.

594 595 596	Ben Hamner, David Fullerton, Kevin Montrose, Rebecca Chernoff, and Will Cole. Predict closed questions on stack overflow, 2012. URL https://kaggle.com/competitions/predict-closed-questions-on-stack-overflow.
597 598 599	Jacob Hansen and Robert Ghrist. Toward a spectral theory of cellular sheaves. <i>Journal of Applied and Computational Topology</i> , 3(3-4):315–358, 2019.
600 601	Zellig S. Harris. Distributional structure. WORD, 10(2-3):146–162, 1954.
602 603 604 605	Lianzhe Huang, Dehong Ma, Sujian Li, Xiaodong Zhang, and Houfeng Wang. Text level graph neu- ral network for text classification. In <i>Proceedings of the 2019 Conference on Empirical Methods</i> <i>in Natural Language Processing and the 9th International Joint Conference on Natural Language</i> <i>Processing (EMNLP-IJCNLP)</i> , pp. 3444–3450. Association for Computational Linguistics, 2019.
606 607 608	Ashish Jaiswal, Ashwin Ramesh Babu, Mohammad Zaki Zadeh, Debapriya Banerjee, and Fillia Makedon. A survey on contrastive self-supervised learning. <i>Technologies</i> , 9(1), 2021.
609 610 611	Yoon Kim. Convolutional neural networks for sentence classification. In <i>Proceedings of the 2014</i> <i>Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pp. 1746–1751. Association for Computational Linguistics, October 2014.
612 613 614	Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional net- works. In <i>International Conference on Learning Representations (ICLR)</i> , 2017.
615 616 617 618	Yuxiao Lin, Yuxian Meng, Xiaofei Sun, Qinghong Han, Kun Kuang, Jiwei Li, and Fei Wu. Bert- GCN: Transductive text classification by combining GNN and BERT. In <i>Findings of the As-</i> <i>sociation for Computational Linguistics: ACL-IJCNLP 2021</i> , pp. 1456–1462. Association for Computational Linguistics, 2021.
619 620 621 622	Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. In <i>International Conference on Learning Representations</i> , 2017.
623 624 625 626 627	Hu Linmei, Tianchi Yang, Chuan Shi, Houye Ji, and Xiaoli Li. Heterogeneous graph attention networks for semi-supervised short text classification. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pp. 4821–4830, Hong Kong, China, November 2019. Association for Computational Linguistics.
628 629 630	Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. Recurrent neural network for text classification with multi-task learning. In <i>Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence</i> , IJCAI'16, pp. 2873–2879. AAAI Press, 2016.
632 633 634	Xien Liu, Xinxin You, Xiao Zhang, Ji Wu, and Ping Lv. Tensor graph convolutional networks for text classification. In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020</i> , pp. 8409–8416. AAAI Press, 2020.
635 636 637 638 639	Yonghao Liu, Renchu Guan, Fausto Giunchiglia, Yanchun Liang, and Xiaoyue Feng. Deep attention diffusion graph neural networks for text classification. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pp. 8142–8152. Association for Computational Linguistics, 2021.
640 641 642	Yonghao Liu, Lan Huang, Fausto Giunchiglia, Xiaoyue Feng, and Renchu Guan. Improved graph contrastive learning for short text classification. <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 38(17):18716–18724, Mar. 2024.
643 644 645 646	Bo Pang and Lillian Lee. A sentimental education: Sentiment analysis using subjectivity summa- rization based on minimum cuts. In <i>Proceedings of the 42nd Annual Meeting on Association for</i> <i>Computational Linguistics</i> , pp. 271–278. Association for Computational Linguistics, 2004.
647	Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In <i>Proceedings of the ACL</i> , 2005.

663

664

682

688

689

690

691

696

697

698

648	Jeffrey Pennington, Richard Socher, and Christopher D. Manning, Glove: Global vectors for word
649	representation. In <i>Empirical Methods in Natural Language Processing (EMNLP)</i> , pp. 1532–1543.
650	2014.
651	

- Kuan-Hieu Phan, Le-Minh Nguyen, and Susumu Horiguchi. Learning to classify short and sparse text & web with hidden topics from large-scale data collections. In *Proceedings of the 17th International Conference on World Wide Web*, WWW '08, pp. 91–100, New York, NY, USA, 2008. Association for Computing Machinery.
- Anand Rajaraman and Jeffrey David Ullman. *Data Mining*, pp. 1–17. Cambridge University Press, 2011.
- Patrick Reiser, Marcel Neubert, Andreas Eberhard, Lorenzo Torresi, Cheng Zhou, Cheng Shao, Hadi
   Metni, Christof van Hoesel, Henning Schopmans, Tobias Sommer, and Pascal Friederich. Graph
   neural networks for materials science and chemistry. *Communications Materials*, 3(1):93, 2022.
  - T. Mitchell Roddenberry, Nicholas Glaze, and Santiago Segarra. Principled simplicial neural networks for trajectory prediction. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pp. 9020–9029. PMLR, 2021.
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- Michael T. Schaub, Jean-Baptiste Seby, Florian Frantzen, Thomas M. Roddenberry, Yu Zhu, and
   Santiago Segarra. Signal processing on simplicial complexes. In *Higher-Order Systems. Under- standing Complex Systems*, pp. 285–309. Springer, Cham, 2022.
- <sup>671</sup>
   <sup>672</sup>
   <sup>673</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>675</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>675</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>674</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>677</sup>
   <sup>678</sup>
   <sup>678</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>671</sup>
   <sup>672</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>677</sup>
   <sup>678</sup>
   <sup>678</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>670</sup>
   <sup>671</sup>
   <sup>672</sup>
   <sup>672</sup>
   <sup>673</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>677</sup>
   <sup>678</sup>
   <sup>678</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>679</sup>
   <sup>671</sup>
   <sup>671</sup>
   <sup>672</sup>
   <sup>672</sup>
   <sup>673</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>674</sup>
   <sup>675</sup>
   <sup>675</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>676</sup>
   <sup>676</sup>
   <sup>677</sup>
   <sup>676</sup>
   <sup>676</sup>
- Freyr Sverrisson, Jean Feydy, Bruno E Correia, and Michael M Bronstein. Fast end-to-end learning
   on protein surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15272–15281, 2021.
- Jian Tang, Meng Qu, and Qiaozhu Mei. Pte: Predictive text embedding through large-scale heterogeneous text networks. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '15, pp. 1165–1174, New York, NY, USA, 2015.
  Association for Computing Machinery.
- Joshua Townsend, Christopher P. Micucci, John H. Hymel, et al. Representation of molecular structures with persistent homology for machine learning applications in chemistry. *Nature Communications*, 11(1):3230, 2020.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
   Bengio. Graph attention networks. *International Conference on Learning Representations*, 2018.
  - Jin Wang, Zhongyuan Wang, Dawei Zhang, and Jun Yan. Combining knowledge with deep convolutional neural networks for short text classification. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pp. 2915–2921, 2017.
- Yaqing Wang, Song Wang, Quanming Yao, and Dejing Dou. Hierarchical heterogeneous graph
   representation learning for short text classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3091–3101, Online and Punta Cana,
   Dominican Republic, November 2021. Association for Computational Linguistics.
  - JunJie Wee and Kelin Xia. Persistent spectral based ensemble learning (PerSpect-EL) for protein-protein binding affinity prediction. *Briefings in Bioinformatics*, 23(2):bbac024, 02 2022.
- Zhihao Wen and Yuan Fang. Augmenting low-resource text classification with graph-grounded pre training and prompting. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '23, pp. 506–516, New York, NY, USA, 2023. Association for Computing Machinery.

702	Kelin Xia and Guo-Wei Wei. Persistent homology analysis of protein structure, flexibility, and
703	folding. International Journal for Numerical Methods in Biomedical Engineering, 30(8):814–
704	844, 2014.
705	

- Jiaming Xu, Bo Xu, Peng Wang, Suncong Zheng, Guanhua Tian, and Jun Zhao. Self-taught con volutional neural networks for short text clustering. *Neural networks : the official journal of the International Neural Network Society*, 88:22–31, 2017.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- Liang Yao, Chengsheng Mao, and Yuan Luo. Graph convolutional networks for text classification.
   In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence*. AAAI Press, 2019.
- Rui Ye, Xin Li, Yujie Fang, Hongyu Zang, and Mingzhong Wang. A vectorized relational graph convolutional network for multi-relational network alignment. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pp. 4135–4141. AAAI Press, 2019.
- Zhihao Ye, Gongyao Jiang, Ye Liu, Zhiyong Li, and Jin Yuan. Document and word representations generated by graph convolutional network and bert for short text classification. In *European Conference on Artificial Intelligence*, 2020.
- Yufeng Zhang, Xueli Yu, Zeyu Cui, Shu Wu, Zhongzhen Wen, and Liang 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*, pp. 334–339. Association
   for Computational Linguistics, 2020.
- Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. A robustly optimized BERT pre-training approach with post-training. In *Proceedings of the 20th Chinese National Conference on Computational Linguistics*, pp. 1218–1227, Huhhot, China, August 2021. Chinese Information Processing Society of China.
  - A. Zia, A. Khamis, J. Nichols, et al. Topological deep learning: a review of an emerging paradigm. *Artificial Intelligence Review*, 57:77, 2024. doi: 10.1007/s10462-024-10710-9.

#### A EFFICIENCY STUDIES

In order to study computational efficiency with the inclusion of higher-order objects, we compute the number of trainable parameters, as shown in Table 4.

	Tuble 1. Humber of Humberb.									
	Twitter	MR	Snippets	StackOverflow						
SCN – 0-simplex	3,654,154	6,120,154	9,446,428	3,969,676						
SCN – 1-simplex	3,671,050	6,137,050	9,463,324	3,986,572						
SCN	3,672,202	6,138,202	9,464,476	3,987,724						
$\overline{C}-\overline{SCN} = \overline{O}-\overline{Simplex}$	113,236,620	115,702,620	119,029,668	113,554,464						
C-SCN – 1-simplex	113,253,516	115,719,516	119,046,564	113,571,360						
C-SCN	113,255,568	115,720,668	119,047,716	113,572,512						

Table 4: Number of trainable parameters.

The time to complete training and evaluation after ten iterations in seconds is also included for analysis, as shown in Table 5.

It is observed that when adding 1-simplexes and 2-simplexes to SCN step-by-step, the average number of trainable parameters increases by 0.18%, and the time increases by 10.55% on average. For
C-SCN, the number of trainable parameters increases by less than 0.1% on average, and the time for
training increases by 3.32% on average. The results demonstrate the computational efficiency of our model involving higher-order complexes in representation learning.

ſ	5	6
7	5	7
7	5	8

Table 5: Time to complete training and evaluation after ten iterations.

	-	0		
	Twitter	MR	Snippets	StackOverflow
SCN – 0-simple	x 639	730	1,009	1,448
SCN – 1-simple	x 679	743	1,235	1,728
SCN	750	798	1,279	1,956
$\overline{C}-\overline{SCN}-\overline{O}-\overline{simp}$	$ e\bar{x} ^{-} \overline{1}, \overline{0}1\overline{5} $	1,033	<u>1,</u> 943	
C-SCN – 1-simp	lex 1,025	1,152	1,970	3,191
C-SCN	1,040	1,197	2,007	3,384

#### **B** Additional Results for Ablation Studies

#### B.1 COMPARED WITH CONTEXTUAL EMBEDDINGS IN CONTRASTIVE LEARNING

Instead of fixed GloVe embedding for word nodes, we compare the results with contextual embeddings (Cont. Emb.) from the BERT model in the following table.

Table 6: Results of test accuracy to compare with the separate contrastive loss

Dataset	Twitter		MR		Snippets		StackOverflow	
Metrics	F1	Acc	F1	Acc	F1	Acc	F1	Acc
BERT	54.92	51.16	51.69	50.65	79.31	78.47	66.94	67.26
C-SCN - Cont. Emb.	74.60	75.01	50.46	55.34	83.31	83.96	82.53	83.00
C-SCN	75.61	76.09	69.46	<b>69.8</b> 7	84.97	85.56	84.15	83.87

It is observed that C-SCN with fixed embeddings achieves better results than the one with contextual embeddings. One explanation could be the limited number of higher-order objects formed with contextual embeddings. 0-simplexes (nodes), which refer to the same word, will not be seen as the same 0-simplex at different locations in the document with contextual embeddings. This will lead to no 2-simplexes formed in the document since one 0-simplex will not be connected again within the text, limiting the expressiveness of structural representations of the higher-order objects.

#### B.2 SEPARATE CONTRASTIVE LOSS FROM THE OBJECTIVE

To evaluate the effectiveness of optimising the contrastive loss and objective function together, experiments to separate the two losses (Sep. Loss) are also conducted. The contrastive loss is first minimised for 100 epochs without labels, and the loss against the final label is optimised with early stopping. The result is shown in the following table.

Table 7: Results of test accuracy to compare with and without GRU.

Dataset	Twitter		MR		Snippets		StackOverflow	
Metrics	F1	Acc	F1	Acc	F1	Acc	F1	Acc
BERT	54.92	51.16	51.69	50.65	79.31	78.47	66.94	67.26
C-SCN - Sep. Loss	67.54	68.35	53.61	56.79	64.69	64.8	27.01	29.78
C-SCN	75.61	76.09	69.46	<b>69.87</b>	84.97	85.56	84.15	83.87

It is observed that with limited training samples (20 samples from each category), pre-training with
 a contrastive loss followed by supervised training does not help the model improve. In detail, the
 separate contrastive loss could improve BERT's performance in binary classification in Twitter and
 MR datasets. In contrast, it worsens the performance in multi-label classification, and the most
 deterioration is from the StackOverflow dataset, which has 20 categories.

## 810 B.3 THE ROLE OF GRU IN THE MESSAGE FUNCTION

We adopted GRU as the UPDATE function to control the amount of information from the previous step and aggregated neighbourhood information. This is achieved through the reset gate and the reset gate structure in GRU. In contrast, we study the role of GRU by comparing the performance if we remove GRU as the UPDATE function and replace it with the sum aggregation (SUM).

Table 8: Results of test accuracy to compare with contextual embeddings.

							U	
Dataset	Twitter		MR		Snippets		StackOverflow	
Metrics	F1	Acc	F1	Acc	F1	Acc	F1	Acc
SCN - SUM	62.3	63.21	54.99	56.23	77.06	77.05	76.02	73.73
SCN	66.13	67.25	61.15	61.93	76.13	75.66	76.85	74.04
$\overline{C}-\overline{S}\overline{C}\overline{N}-\overline{S}\overline{U}\overline{M}$	74.01	74.45	55.51	58.26	77.38	77.92	78.73	77.46
C-SCN	75.61	76.09	69.46	<b>69.87</b>	84.97	85.56	84.15	83.87

One challenge we observed without GRU was the overfitting issue on the train set across different datasets. The results deteriorated when we removed GRU from SCN and C-SCN respectively, illustrating the importance of GRU in the message-passing mechanism for higher-order complexes.

- B.4 PSEUDO-CODE FOR C-SCN
- We include the pseudo-code for C-SCN to enhance the reproducibility.

Algorithm 1: Algorithm Pseudo Code for C-SCN.

**Input:** Text data with *words, punctuations and label* as shown in Figure 3.

- <sup>837</sup> Tokenised words from the document data  $\mathbf{x}_{doc}$ ;
- Tokenised unique 0-simplex  $\sigma_0 \in \mathcal{K}_0$ ;
- 1-simplex indices following the chronological order of tokens  $\sigma_1 \in \mathcal{K}_1$ ;

1-simplex features tokenised to one of the types: *forward, backward, self-loop*;

- 2-simplex features tokenised by the flow directions of components;
- 842 /\* Add higher-order simplexes if needed.

\*/

## 843 <u>Model Construction</u>

844 Parameters:

816 817

826

827

828 829

830

833

834

- Embedding matrices for 0-simplexes, 1-simplexes and 2-simplexes:  $\mathcal{E}_0$ ,  $\mathcal{E}_1$  and  $\mathcal{E}_2$ ;
- 846 Number of layers: L;
- Message-passing mechanism for 0-simplexes and 1-simplexes following Equation 2 and Equation 3:  $MP_0$ ,  $MP_1$ ;
- Attention mechanism for 0-simplexes and 1-simplexes following Equation 4: Attn<sub>0</sub>, Attn<sub>1</sub>; Activation function:  $\phi$ ;
- 850 Activation function:  $\phi$ , Transformer model: Trans<sub>t</sub>;
- Linear layers that process the output of the transformer and the SCN to the label space:  $MLP_t$ , MLP<sub>s</sub>.

#### <sup>853</sup> Initialise features:

 $\mathbf{h}_{\sigma_0}^{(0)} = \mathcal{E}_0(\sigma_0) \,\forall \, \sigma_0 \in \mathcal{K}_0; \, \mathbf{h}_{\sigma_1}^{(0)} = \mathcal{E}_1(\sigma_1) \,\forall \, \sigma_1 \in \mathcal{K}_1; \, \mathbf{h}_{\sigma_2}^{(0)} = \mathcal{E}_2(\sigma_2) \,\forall \, \sigma_2 \in \mathcal{K}_1.$ 854 855 for  $\ell = 1$  to L - 1 do  $\begin{bmatrix} \mathbf{h}_{\sigma_1}^{(\ell)} = \mathrm{MP}_1(\mathbf{h}_{\sigma_1}^{(\ell-1)}, \mathbf{h}_{\sigma_2}^{(\ell-1)}); \\ \mathbf{h}_{\sigma_0}^{(\ell)} = \mathrm{MP}_0(\mathbf{h}_{\sigma_0}^{(\ell-1)}, \mathbf{h}_{\sigma_1}^{(\ell-1)}). \end{bmatrix}$ 856 857 858  $\mathbf{h}_{\mathcal{K}}^{(L)} = \sum_{\sigma_1 \in \mathcal{K}_1} \text{Attn}_1(\mathbf{h}_{\sigma_1}^{(L)}) \oplus \sum_{\sigma_0 \in \mathcal{K}_0} \text{Attn}_0(\mathbf{h}_{\sigma_0}^{(L)});$ 859  $Z_t = \text{MLP}_t(\text{Trans}_t(\mathbf{x}_{\text{doc}})); \ Z_s = \text{MLP}_s(\mathbf{h}_{\mathcal{K}}^{(L)});$ 860 861  $Z = \frac{1}{2} \left( \operatorname{softmax}(Z_s) + \operatorname{softmax}(Z_t) \right);$ 862 return  $\hat{y} = Z, Z_s, Z_t$ .

<sup>836</sup> Simplicial Complex Construction