

MIST: Mutual Information Maximization for Short Text Clustering

Anonymous ACL submission

Abstract

Short text clustering poses substantial challenges due to the limited amount of information provided by each text sample. Previous efforts based on dense representations are still inadequate as texts are not sufficiently segregated in the embedding space before clustering. Even though the state-of-the-art method utilizes contrastive learning to boost performance, the process of summarizing all local tokens to form a sequence representation for the whole text includes noise that may obscure limited key information. We propose Mutual Information Maximization Framework for Short Text Clustering (MIST), which overcomes the information drown-out by including a mechanism to maximize the mutual information between representations on both sequence and token levels. Experimental results across eight standard short text datasets show that MIST outperforms the state-of-the-art method in terms of Accuracy or Normalized Mutual Information in most cases.

1 Introduction

Text clustering is a vital task for a wide range of downstream applications. It aims to partition texts into groups of similar categories in an unsupervised manner. The growth of social media, discussion forums, and news aggregator websites has led to a large number of short-length texts being produced daily. Therefore, clustering these short texts has become crucial for many real-world applications ranging from recommendation to text retrieval (Yohannes and Assabie, 2021).

In short texts, the most informative words and phrases of the text content usually appear only once. This exacerbates the sparsity problem, posing an additional hurdle for clustering short texts. Traditional methods, such as BoW and TF-IDF, provide relatively sparse representation vectors with limited descriptive power. Hence, they perform poorly when clustered using a standard distance-based clustering algorithm (Hadifar et al., 2019).

To address this problem, most recent methods (Xu et al., 2017; Hadifar et al., 2019; Yin et al., 2021) utilize deep neural networks to map high-dimensional data into meaningful dense representations in a lower-dimensional space and adopt a multi-stage scheme in which the clustering process is performed after learning feature representations. However, the clustering performance of these methods remains unsatisfactory as texts still have a lot of overlap among categories in the latent space before clustering (Zhang et al., 2021).

Alternatively, an end-to-end clustering scheme (Zhang et al., 2021; Xie et al., 2016) simultaneously optimizes representation learning and clustering objectives. To achieve desirable outcomes, Zhang et al. (2021) propose a method that employs contrastive representation learning, which has been successful in self-supervised learning and can help spread out overlapping categories, in order to obtain effective short text representations.

As shown in Zhang et al. (2021), improving representation is crucial for enhancing the clustering performance. Nevertheless, the contrastive learning method used in Zhang et al. (2021) only considers sequence-level embeddings that are formed by averaging all local tokens in each text instance, including uninformative noise. This could generate a representation in which sparse yet informative terms used to describe the text content may be obscured by noise, potentially affecting the clustering performance. We consider the preservation of limited information in such a low signal-to-noise environment as a vital feature for short-text clustering. Addressing this gap will result in sequence representations that are more semantically representative and robust to noisy tokens in short texts.

In this paper, we introduce the Mutual Information Maximization Framework for Short Text Clustering (MIST), a new multi-stage approach. We aim to improve representation learning stage for short text clustering using two contrastive

083 learning objectives operating at the sequence and
084 token levels. In particular, we apply the concept of
085 mutual information (MI) maximization to facilitate
086 us in comparing the semantic similarity between
087 representations across the two hierarchical levels.

088 The crux of our method lies in integrating the
089 *sequence-level* and *token-level* MI maximization
090 objectives concurrently for the following purposes.

- 091 1. *Learning Distinct Text Representation*: The
092 first learning objective maximizes MI between
093 each positive pair at the sequence level;
- 094 2. *Informative Token Preservation*: The second
095 objective is designed to enforce each text rep-
096 resentation at the sequence level to extract lo-
097 cal information shared across all its individual
098 tokens by directly maximizing MI between
099 them. This way, we mitigate the obscurity-
100 by-noise problem and preserve limited key
101 information in a weak signal environment.

102 The growth in the size of short text sequences
103 may exacerbate a poor signal-to-noise ratio. To
104 deal with short text samples with various signal-to-
105 noise ratios, we additionally propose an *adaptive*
106 *weighting function* that dynamically determines an
107 appropriate ratio between the two objectives based
108 on the length of the texts. To our knowledge, the
109 method of combining two MI maximization objec-
110 tives logically is presented for the first time. Note
111 that the representations at different levels have a
112 direct implication on one another, and the sequence
113 representations are subsequently used in the clus-
114 tering stage by applying the *k*-means algorithm.

115 We conduct extensive experimental studies over
116 the eight standard benchmarks. MIST improves
117 the clustering performance in terms of Accuracy
118 and Normalized Mutual Information in most cases
119 compared to the current state-of-the-art while using
120 an identical configuration across all datasets. This
121 demonstrates the generalizability of our method.

122 Our main contributions are outlined as follows:
123 (1) We propose a novel representation learning tech-
124 nique for short text clustering through the integra-
125 tion of sequence-level and token-level MI maxi-
126 mization objectives. (2) To balance the two objec-
127 tives, we introduce an adaptive weighting func-
128 tion. (3) Our ablation study provides a further
129 demonstration of how different prioritization of
130 the two MI objectives impacts the clustering per-
131 formance across datasets of various text lengths;
132 as text length increases, the preservation of limited
133 local information becomes more significant.

2 Related Work 134

Short Text Clustering. There are several strate- 135
gies to overcome the sparsity of short text represen- 136
tations. Some recent methods utilize a multi-stage 137
architecture that breaks down the clustering frame- 138
work into multiple stages; the clustering process 139
is performed after learning feature representations. 140
Xu et al. (2015, 2017) use a convolutional neural 141
network to learn non-biased representations by fit- 142
ting the output units with pretrained-binary codes 143
from a dimensionality reduction method. Hadi- 144
far et al. (2019) utilize Smooth Inverse Frequency 145
(Arora et al., 2017) to obtain weighted word embed- 146
dings. During training, they enrich discriminative 147
features by tuning an autoencoder with soft clus- 148
tering assignments. For the aforementioned works, 149
the *k*-means clustering is employed on the trained 150
representations to get the final clusters. 151

152 Another approach is to enhance the quality of
153 the initial clustering with an iterative classification
154 algorithm. Rakib et al. (2020) proposed the ECIC
155 algorithm to detect and remove outliers in each
156 iteration. Moreover, they make use of word em-
157 beddings by averaging them to represent each text,
158 and combine the ECIC algorithm with hierarchical
159 clustering. To boost the clustering quality further,
160 Pugachev and Burtsev (2021) exploit deep sentence
161 representations (Cer et al., 2018) and make modifi-
162 cations to the ECIC algorithm.

163 The recent state-of-the-art, SCCL (Zhang et al.,
164 2021), leverages contrastive learning to encourage
165 greater separation between overlapped categories
166 in the original data space. By jointly optimizing a
167 contrastive loss and a clustering objective (Reimers
168 and Gurevych, 2019a), SCCL outperforms prior
169 works and yields cutting-edge results. In addition,
170 other constrastive learning methods have also been
171 experimented on short text clustering, such as using
172 entities for contrastive learning to provide supervi-
173 sion signals for their related sentences (Nishikawa
174 et al., 2022), and using virtual augmentation for
175 contrastive learning to circumvent the discrete na-
176 ture of language (Zhang et al., 2022).

177 Moreover, a new technique for short text clus-
178 tering is presented in Zheng et al. (2023); it comprises
179 a pseudo-label generation module and a robust rep-
180 resentation learning module. The former generates
181 pseudo-labels, which are robust against the imbal-
182 ance in data, as the supervision for the latter.

Self-Supervised Learning. Self-supervision has 183
gained popularity and become a common technique 184

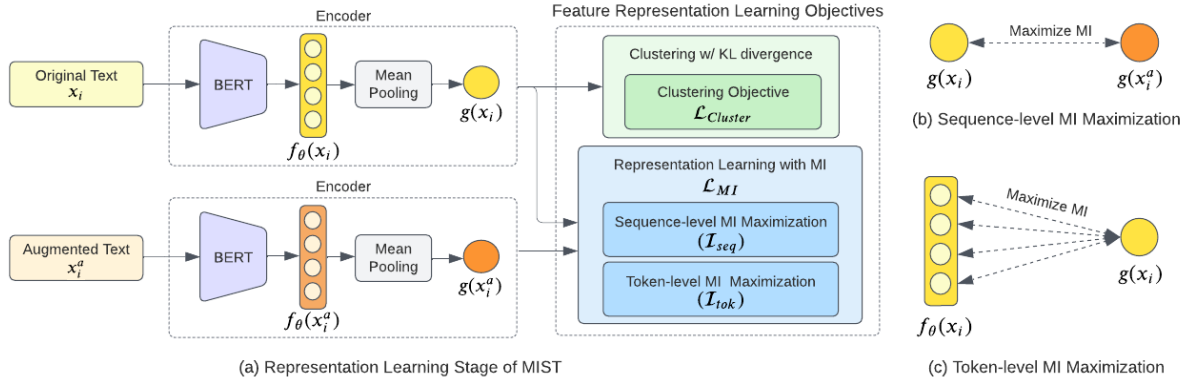


Figure 1: (a) Representation Learning Stage Overview. MIST considers all pairs of original text x_i , and its augmented version x_i^a as positive samples. MIST jointly optimizes the clustering objective $\mathcal{L}_{Cluster}$, and the MI objective \mathcal{L}_{MI} , which includes (b) a sequence-level MI maximization objective \mathcal{I}_{seq} that maximizes MI between representations at the sequence level (x_i and x_i^a), and (c) a token-level MI maximization objective \mathcal{I}_{tok} that directly maximizes MI between a sequence representation (of both x_i and x_i^a) and its tokens ($f_\theta(x_i)$ and $f_\theta(x_i^a)$).

in unsupervised representation learning for a variety of downstream purposes (Chen et al., 2020; He et al., 2020; Caron et al., 2020; Grill et al., 2020).

Learning meaningful representations by estimating and maximizing MI is one of the prominent contrastive learning strategies. Its effectiveness has been demonstrated in both vision (Hjelm et al., 2019; Bachman et al., 2019; Sordoni et al., 2021) and text domains (Kong et al., 2020; Caron et al., 2020; Giorgi et al., 2021). Deep Infomax (DIM) (Hjelm et al., 2019) introduces global and local MI maximization objectives for learning image representations. However, each of these is implemented separately according to the task. The authors find success in optimizing local MI maximization objective by maximizing MI between local features and global features. Inspired by local Deep InfoMax, Zhang et al. (2020) proposes a sentence representation learning method that maximizes the MI between the sentence-level representation and its CNN-based n-gram contextual dependencies.

On the contrary, we integrate two MI maximization strategies concurrently to learn textual representations for various short text characteristics. We also introduce a generalized adaptive weighting function for effectively integrating both objectives.

3 Proposed Method: MIST

We propose a short text clustering framework consisting of two stages. First, we train a model using feature representation learning objectives as illustrated in Figure 1. Second, we apply the k -means algorithm on the trained representations at inference time to obtain the final clusters. This

investigation focuses on improving the first stage.

The main idea of our solution lies in the learning objective function \mathcal{L} that takes into account an MI objective \mathcal{L}_{MI} and an unsupervised clustering objective $\mathcal{L}_{Cluster}$, which is used to enforce the encoder to capture categorical structure and provide a suitable representation space for clustering task.

$$\mathcal{L} = \beta \mathcal{L}_{MI} + \eta \mathcal{L}_{Cluster}, \quad (1)$$

where β and η represent the trade-off between \mathcal{L}_{MI} , and $\mathcal{L}_{Cluster}$. In our experiments, we set β to 1 and η to 2 to provide more weight to $\mathcal{L}_{Cluster}$.

In Section 3.1, we describe our main contribution, the MI maximization learning procedure, including (1) sequence-level and token-level MI maximization objectives; (2) an adaptive weighting function that is also incorporated to balance them. Section 3.2 presents the auxiliary clustering objective utilized in the learning stage.

3.1 Representation Learning with MI Maximization

Short texts are challenging to cluster due to the weak signal caused by noise. In the context of this study, short texts are recognized as those that are short in length and typically contain informal fragmental non-sentence structures, e.g., tweets and news snippets. One strategy to improve the clustering performance is to adopt *contrastive learning* to construct an embedding space that minimizes local invariance for each positive pair. However, a standard contrastive learning procedure, which is performed by contrasting between sequence representations (global features), may allow noise to *drown out* sparse but informative local-token embeddings (local features) when these tokens are

mean-pooled to form a sequence representation. Consequently, optimizing solely contrastive learning at the sequence level is insufficient for learning representations in a weak signal environment.

3.1.1 Hierarchical MI Objective

In contrast to previous works on MI maximization learning, which utilized each MI objective separately, we incorporate the learning of both sequence and token representations into a single objective. This strategy offers two advantages: (1) it mitigates the problem of information drown-out by allowing individual tokens to participate in the MI maximization process; (2) it supports weight adjustment between these two MI levels to handle short text inputs with various signal-to-noise ratios.

Sequence-Token MI Maximization. According to Tian et al. (2020), contrastive learning is equivalent to maximizing the lower bound of MI between a sequence representation and its augmented version (positive). Intuitively, it reflects how much more precisely we can determine the representation given a positive compared to when we are unaware of the positive (Bachman et al., 2019). This principle enables us to incorporate an additional mechanism beyond the sequence-level objective.

We build our framework based on the MI maximization concept through the integration of two MI objectives. In this way, our model can effectively learn distinct short text representations using the *sequence-level* MI objective while simultaneously preserving local information using an additional objective. Specifically, the *token-level* MI objective helps alleviate the information obscurity from noise by maximizing the MI between each local token and its sequence representation. As a result, the overall learning objective \mathcal{L}_{MI} consists of two components: (1) sequence-level MI maximization \mathcal{I}_{seq} , and (2) token-level MI maximization \mathcal{I}_{tok} , operating concurrently in a sequence-token hierarchy as shown in Figure 1.

$$\mathcal{L}_{\text{MI}} = -(1 - \lambda)\mathcal{I}_{\text{seq}} - \lambda\mathcal{I}_{\text{tok}}, \quad (2)$$

where λ corresponds to the balancing weight for \mathcal{I}_{seq} and \mathcal{I}_{tok} objectives, which is defined in Eq.3.

Adaptive Weighting Function. According to our analysis, short text inputs vary in length across different datasets, ranging from *fragmental sequences* of 6 words to 28 words. Regarding signal-to-noise, larger sequences tend to contain a greater proportion of noise that does not provide useful semantics for the clustering step, whereas informative

terms *usually still appear once*. This exacerbates the information drown-out problem due to a poor signal-to-noise ratio.

While other short text clustering techniques treat all text samples in the same fashion, we argue that different-length short texts should be handled differently. We propose an MI maximization strategy adaptable to text length so that our method can efficiently deal with short text instances containing varying signal-to-noise ratios, without the need for a hyperparameter search for any particular dataset. Since larger sequences necessitate more effort to preserve limited crucial information, we place more weight on the \mathcal{I}_{tok} objective by encouraging λ to be larger as the total number of tokens in the text grows. Thus, our *generalized adaptive weighting function* (Eq.3) is introduced to assign the weight of λ depending on the average number of tokens in text samples for each minibatch of size N :

$$\lambda = \max\left(0, \left[\frac{0.1}{N} \sum_{i=1}^N l_i\right] - 1\right) \times 0.1, \quad (3)$$

where l_i denotes the number of tokens in a text x_i and it is directly proportional to the text length.

In the representation learning stage, we randomly sample a minibatch $X^o = x_1^o, \dots, x_N^o$ of N original texts with empirical probability distribution \mathbb{P} . Then, we generate an augmented version for each text to obtain an augmented batch $X^a = x_1^a, \dots, x_N^a$, where X^o and X^a are of identical size. The encoder f_θ , a pretrained transformer network, encodes an input text x into a sequence of contextualized token embeddings with length l , $f_\theta(x) := \{f_\theta^{(i)}(x) \in \mathbb{R}^d\}_{i=1}^l$, where i is the token index and d is the number of dimension. These token representations are then mean pooled $m(f_\theta(x))$ to generate a sequence representation denoted as $g(x) = m(f_\theta(x)) \in \mathbb{R}^d$.

3.1.2 Computing the Sequence-level MI.

This learning objective, \mathcal{I}_{seq} , aims to learn distinct text representations through the maximization of MI between the original sample and its augmented version at the sequence level. By treating each original text $g(x^o)$ and its augmentation $g(x^a)$ as positive pairs, the \mathcal{I}_{seq} objective is defined as:

$$\mathcal{I}_{\text{seq}} = \frac{1}{N} \left(\sum_{x \in X} \hat{\mathcal{I}}^{JSD}(g(x^o); g(x^a)) \right) \quad (4)$$

We adopt a Jensen-Shannon estimator (Nowozin et al., 2016; Hjelm et al., 2019) to estimate a lower bound of MI, $\hat{\mathcal{I}}_\theta^{JSD}$:

$$\begin{aligned} \widehat{\mathcal{I}}_{\theta}^{JSD}(g(x^o); g(x^a)) := & \\ E_{\mathbb{P}}[-sp(-g(x^o) \cdot g(x^a))] & \quad (5) \\ - E_{\mathbb{P} \times \tilde{\mathbb{P}}} [sp(g(x^o) \cdot g(\tilde{x}^a))], & \end{aligned}$$

where \tilde{x}^a is a negative augmented textual input sampled from distribution $\tilde{\mathbb{P}} = \mathbb{P}$, and $sp(z) = \log(1 + e^z)$ is the softplus function.

3.1.3 Computing the Token-level MI

In contrast to Zhang et al. (2020), we constrain the sequence representation containing high MI with each token to preserve limited local information in short texts— by maximizing MI between the sequence representation and its token representations directly— instead of its local contextual n-gram embeddings. In particular, we attempt to maximize the average MI between a sequence representation and all its token representations while minimizing MI with the tokens of other texts. We define \mathcal{I}_{tok} for each minibatch as

$$\begin{aligned} \mathcal{I}_{tok} = \frac{1}{2N} \left(\sum_{x^o \in X^o} \sum_{i=1}^{l_{x^o}} \widehat{\mathcal{I}}^{JSD}(g(x^o); f_{\theta}^{(i)}(x^o)) \right. \\ \left. + \sum_{x^a \in X^a} \sum_{i=1}^{l_{x^a}} \widehat{\mathcal{I}}^{JSD}(g(x^a); f_{\theta}^{(i)}(x^a)) \right). \end{aligned} \quad (6)$$

An estimated MI for each sequence $g(x)$ and token representations $f_{\theta}^{(i)}(x)$ is calculated as follows:

$$\begin{aligned} \widehat{\mathcal{I}}_{\theta}^{JSD}(g(x); f_{\theta}^{(i)}(x)) := & \\ E_{\mathbb{P}}[-sp(-g(x) \cdot f_{\theta}^{(i)}(x))] & \quad (7) \\ - E_{\mathbb{P} \times \tilde{\mathbb{P}}} [sp(g(x) \cdot f_{\theta}^{(i)}(\tilde{x}))], & \end{aligned}$$

where \tilde{x} is a different text on the minibatch.

3.2 Clustering with KL Divergence

In addition to the MI objective, we employ $\mathcal{L}_{\text{Cluster}}$ during the learning stage to encourage the coalescence of samples that are most likely to belong to the same cluster. We follow the clustering method proposed by Xie et al. (2016), which is also used by Zhang et al. (2021). This method involves computing soft cluster assignments and formulating the clustering objective using KL divergence.

For the first step, we follow Xie et al. (2016) using the Student’s t-distribution Q to compute a soft cluster assignment for each text instance $x_j \in X$ and the centroid μ_k where $\mu_k \in \{1, \dots, K\}$ for the dataset with K -clusters. Specifically, we compute the probability q_{jk} of assigning a text x_j to a cluster μ_k as follows.

$$q_{jk} = \frac{(1 + \|g(x_j) - \mu_k\|_2^2 / \alpha)^{-\frac{\alpha+1}{2}}}{\sum_{k'=1}^K (1 + \|g(x_j) - \mu_{k'}\|_2^2 / \alpha)^{-\frac{\alpha+1}{2}}} \quad (8)$$

The α symbol represents the degree of freedom of the distribution, and we set α to 1. Following Zhang et al. (2021), each centroid μ_k is approximated by the linear clustering head c_{θ} .

The second step is calculating an auxiliary target distribution P and using it to assist in refining clusters’ centroids. The main idea is to give more importance to text samples with high clustering confidence. The probability $p_{jk} \in P$ is defined as

$$p_{jk} = \frac{q_{jk}^2 / \sum_{j'} q_{j'k}}{\sum_{k'} (q_{jk'}^2 / \sum_{j'} q_{j'k'})}. \quad (9)$$

To match the soft cluster assignments to the target distribution, the KL-divergence between probability distributions P and Q is computed as follows.

$$\ell_j^C = KL[p_j || q_j] = \sum_{k=1}^K p_{jk} \log \frac{p_{jk}}{q_{jk}} \quad (10)$$

We then formulate it as a clustering objective for each minibatch of size N as

$$\mathcal{L}_{\text{Cluster}} = \sum_{j=1}^N \ell_j^C / N. \quad (11)$$

4 Experimental Setup

Datasets. Following previous works (Rakib et al., 2020; Zhang et al., 2021; Pugachev and Burtsev, 2021; Zheng et al., 2023), we conduct experiments and evaluate the performance of MIST on the eight standard short text clustering datasets. The descriptions of the datasets are provided in Appendix A.1

Implementation. We implement our model in PyTorch (Paszke et al., 2017) and use the *paraphrase-mpnet-base-v2* in Sentence Transformers library (Reimers and Gurevych, 2019b) as the encoder, with a linear clustering head following Zhang et al. (2021). The encoder is trained for 1,200 iterations for all datasets and we use Adam optimizer with a batch size of 256. The learning rate of the encoder and the clustering head are set to $6e-6$ and $6e-5$, respectively. We follow Xu et al. (2017) and Hadifar et al. (2019) by randomly selecting 10% of data as the validation set. Furthermore, we follow Zhang et al. (2021) by not performing any preprocessing operations on all eight datasets. Although some of the existing works preprocess the texts by removing symbols, stop words, and punctuation, or converting them to lowercase.

In the training stage, the original and augmented texts are taken into consideration as inputs for the MI objective \mathcal{L}_{MI} , since we found that they are more effective than employing two augmented

	AgNews		SearchSnippets		StackOverflow		Biomedical	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
<i>Reported in the References</i>								
STCC	-	-	77.09	63.16	51.13	49.03	43.62	38.05
Self-Train	-	-	77.1	56.7	59.8	54.8	54.8	47.1
SCA-AE	68.36	34.14	68.71	50.26	76.55	65.99	40.25	33.29
HAC-SD	81.84	54.57	82.69	63.76	64.80	59.48	40.13	33.51
RSTC	84.24	62.45	80.10	69.74	83.30	74.11	48.40	40.12
<i>Reimplementation</i>								
SBERT (k-means)	83.44	57.76	73.02	59.77	76.79	75.12	41.30	36.93
SCCL	85.67	65.98	78.73	70.10	78.35	75.6	39.35	39.2
SCCL-Multi	85.6	66	78.6	70.17	78.3	76.22	39.2	33.7
<i>Proposed Method</i>								
MIST	89.47*	70.25*	76.72	67.69	79.65	78.59*	39.15	34.66
	Tweet		GoogleNews-TS		GoogleNews-T		GoogleNews-S	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
<i>Reported in the references</i>								
STCC	-	-	-	-	-	-	-	-
Self-Train	-	-	-	-	-	-	-	-
SCA-AE	84.85	89.19	-	-	-	-	-	-
HAC-SD	89.62	85.20	85.76	88.00	81.75	84.20	80.63	83.50
RSTC	75.20	87.35	83.27	93.15	72.27	87.39	79.32	89.40
<i>Reimplementation</i>								
SBERT (k-means)	62.7	86.8	67.40	90.47	63.98	86.13	65.87	87.64
SCCL	68.3	88.59	78.9	92.92	69.9	87.9	73.55	89.33
SCCL-Multi	67.55	88.41	80.15	93.4	72.85	88.44	74.2	89.47
<i>Proposed Method</i>								
MIST	91.75*	95.12*	90.63*	96.42*	78.8	89.31*	82.14*	90.86*

Table 1: Experimental results on eight short text clustering datasets. * denotes that MIST is significantly better than both reimplemented versions of SCCL. In order to statistically compare models, we use the Almost Stochastic Dominance test (Dror et al., 2019) with the significant level of 0.05.

433 pairs. We follow Zhang et al. (2021) and utilize
434 *Contextual Augmenter* (Kobayashi, 2018; Ma,
435 2019) to generate augmented samples for each text
436 instance as it was demonstrated to produce the
437 best outcomes in their study. To assess cluster-
438 ing performance, we use the same standard met-
439 rics—Accuracy (ACC) and Normalized Mutual In-
440 formation (NMI)—as used in all competitive meth-
441 ods¹. The results are averaged over five trials.

442 5 Experimental Results

443 We extensively compare the performance of MIST
444 with state-of-the-art methods including STCC (Xu
445 et al., 2017), Self-Train (Hadifar et al., 2019), HAC-
446 SD (Rakib et al., 2020), SCA-AE (Yin et al., 2021),
447 SCCL (Zhang et al., 2021), and RSTC (Zheng et al.,
448 2023). In addition, this section provides an ablation
449 study on our proposed method.

450 5.1 Main Results

451 As shown in Table 1, MIST achieves state-of-the-
452 art results in terms of Accuracy and NMI for most

¹The Accuracy is calculated via the Hungarian algorithm, and NMI measures the information shared between the ground truth assignments and the predicted assignments.

453 cases on the eight benchmark datasets. In contrast,
454 HAC-SD and Self-Train attain the best results in
455 only two cases, whereas SCCL and RSTC produce
456 the best outcome in only one case. Note that, the
457 performances of MIST are collected using the iden-
458 tical setting and training iteration across all datasets
459 to demonstrate generalizability. As a result, the
460 need for a specific configuration for each dataset is
461 avoided, enabling a reduction in model overhead.

462 For datasets with a small number of clusters in
463 the upper section of the table, MIST shows supe-
464 rior performances on AgNews for both metrics and
465 StackOverflow in terms of NMI. Notably, there are
466 two datasets that MIST is outperformed by com-
467 petitors for both ACC and NMI, i.e., Biomedical
468 and SearchSnippets. For Biomedical, Hadifar et al.
469 (2019) dominates the competitive methods. They
470 achieve the best results by using an in-domain pre-
471 trained model to process this dataset, whereas the
472 dataset used to pretrain our encoder and other re-
473 cent methods is a general-domain one.

474 For SearchSnippets, we observe that most of the
475 text samples are collections of keywords and termi-
476 nologies rather than coherent sentence structures.
477 Moreover, SearchSnippets samples are medium-

length fragmental sequences; as a result, the token-level MI maximization objective is more emphasized due to the length of the texts. These two factors exert a direct impact on the token-level MI maximization objective while it is being executed in the learning stage. Since the token vectors are contextualized representations, forcing the model to learn from incoherent contextual signals can be detrimental to the overall sequence representations, which are subsequently used in the clustering stage. This can be more problematic when the same keywords appear in multiple clusters.

As demonstrated in the lower section of the table, MIST obtains the best outcomes on most of the datasets containing a large number of clusters. Due to the fine-grained categorization of these datasets, texts in different clusters may share similar content or keywords, hence inducing ambiguity. This ambiguity in textual data and ground truths leads to erroneous predictions. Moreover, another cause of inaccuracy is when the text content in one cluster is a subtopic of another. GoogleNews-T, which only contains news headlines that are relatively short with few keywords, presents an additional challenge for clustering these extremely short texts into a large number of clusters. In terms of Accuracy, our method achieves a result comparable to that of Rakib et al. (2020) on GoogleNews-T. We conjecture that hierarchical clustering and outlier removal algorithms employed in their method can better deal with the hierarchical nature of data in this scenario. However, MIST outperforms Rakib et al. (2020) in terms of NMI on this dataset.

Although GoogleNews-S and GoogleNews-TS share the same challenges as GoogleNews-T, clustering texts in both datasets is more accurate due to the benefit of additional context and information in the texts themselves. As GoogleNews-S contains snippets of news, and GoogleNews-TS includes both titles and snippets. Consequently, MIST achieves superior clustering performances on both datasets for both matrices.

Furthermore, we thoroughly compare MIST with SCCL, as this current state-of-the-art model also utilizes contrastive learning and aims to improve the effectiveness of representations for short text clustering, which is similar to our contribution, by reproducing SCCL in two versions for a fair comparison: an end-to-end (original) version, and a multi-stage version. For the latter, we apply the k -means algorithm on the trained representations to

get the final clusters, referred to as *SCCL-Multi*. In particular, SCCL-Multi is analogous to our framework, except for the representation learning technique. The reimplemented versions use the same backbone and augmentation setting as our model.

The comparative results show that MIST outperforms SCCL for both versions in most cases. More specifically, the superior performances of MIST compared to SCCL-Multi demonstrate that our proposed representation learning procedure improves short text representations more effectively than the standard contrastive learning objective in the SCCL framework. MIST also consistently surpasses both reimplemented versions of SCCL in other settings, including settings indicated in their publication in most cases, as shown in Appendix A.6.

5.2 Ablation Study

To better understand the effect of the various model modifications on the clustering performance and the analysis versus text lengths, we conducted additional experiments by varying the trade-off between components in our training procedure.

5.2.1 The Impact of Sequence- and Token-MI Maximization Objectives

This experiment studies the impact of the ratio between two MI maximization objectives on the clustering performance and the importance of incorporating both objectives in our representation learning procedure. We report and analyze the performance of our model using four different values of λ . Particularly, λ denotes the weight of token-level MI maximization objectives, \mathcal{I}_{tok} , and $1-\lambda$ represents the weight of sequence-level MI objectives, \mathcal{I}_{seq} . We consider the following settings: (1) *MIST-seq*: our model with a sequence-only MI maximization objective ($\lambda = 0$), (2) *MIST-tok*: our model with a token-only MI maximization objective ($\lambda = 1$), (3) *MIST-equal*: our model with both objectives are given an equivalent weight ($\lambda = 0.5$), and (4) *MIST*: our proposed version, i.e., our model with the λ determined by the adaptive weighting function, Eq.3, varying according to input text length.

As shown in Figure 2, MIST with the value of λ set by Eq.3 yields the best performances in terms of Accuracy for most datasets and shows performance gains compared to other settings. We also discovered that NMI tends to follow the same trend as Accuracy, as presented in Appendix A.2. This demonstrates that the length of short texts has a great impact on determining the appropriate ratio

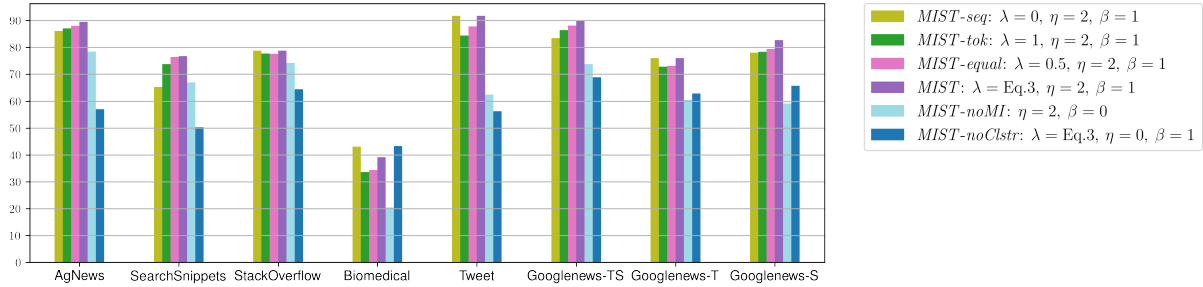


Figure 2: Accuracy for six different settings including four different weighting ratios between sequence-level and token-level MI maximization objectives. As well as, a setting where the clustering objective is absent ($\eta = 0$), and a setting where the MI objective is absent ($\beta = 0$). Note that when we set β to 0, λ has no effect.

between the two MI objectives, i.e. the optimal ratio varies by input samples. By utilizing the proposed adaptive weighting function, MIST can perform effectively across various datasets.

For medium or large fragmental sequences, such as GoogleNews-TS, MIST produces the best outcomes when the weight λ calculated by Eq.3—the value of λ is greater than 0. Remarkably, MIST-equal and MIST-tok always outperform MIST-seq in this situation. This shows that only the sequence-level objective is inadequate when dealing with lengthy texts, as larger fragments usually have a higher signal-to-noise ratio. However, this issue can be mitigated by performing the token-level MI maximization during the learning stage.

For small fragment datasets, such as Tweet, text samples are relatively short and contain less signal-to-noise problem. In this scenario, the weight λ is equal to 0 based on Eq.3, i.e., MIST is identical to MIST-seq, which outperforms all other settings. MIST-tok and MIST-equal may encourage the encoder to learn text representations by placing emphasis on keywords that could also appear in multiple clusters, causing ambiguity and error in clustering. Hence, the token MI objective provides advantages when used in a suitable weight.

In addition, we investigate the situation in which the MI objective is removed ($\beta = 0$), *MIST-noMI*. The ablation results show significant drops in the performance on all datasets. This implies that the MI objective is essential for performance gain.

5.2.2 The Impact of the Clustering Objective

As shown in Figure 2, the clustering performance drops drastically when we remove the clustering objective ($\eta = 0$) during learning representations, *MIST-noClstr*. This demonstrates that the categorical structure imposed by jointly optimizing the clustering objective with the MI objective is a crucial component that boosts performance. Furthermore,

we observe that as the weight of the clustering objective (η) increases, the performances continuously improve until η reaches its saturation point at 2. In Figure 3, the average Accuracy and NMI for all eight datasets improve as the clustering weight is steadily increased until it reaches 2.

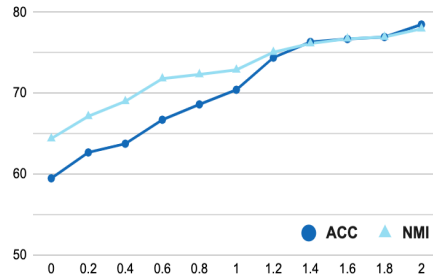


Figure 3: The average clustering performance across eight datasets based on the clustering objective strength.

6 Conclusion

We propose a novel multi-stage short text clustering framework that mainly focuses on improving the representation learning stage. Our adaptive learning approach integrates two MI maximization objectives operating at the sequence and token levels to produce effective representations. This mechanism allows us to simultaneously learn distinct text representations while maintaining limited information in a weak signal environment. In addition, we introduce a generalized adaptive weighting function that considers the length of the texts to determine an optimal ratio between the two MI maximization objectives during the learning stage.

MIST outperforms competitive methods in most cases in terms of Accuracy and NMI across eight benchmark datasets. This demonstrates that utilizing the MI maximization strategy for learning representation in a constrained environment could potentially be a promising tactic. Further study would be worthwhile since it might enhance the quality of textual representations for other tasks.

646 Limitations

647 This section discusses the limitations of our pro-
648 posed framework. Firstly, the encoder of our model
649 is pretrained using general domain data. Hence,
650 the performance drops when our model encounters
651 short texts in a specific domain, such as Biomedical.
652 Furthermore, short text inputs containing only of
653 keywords or incoherent text sequences hinder the
654 performance of our representation learning method.
655 In particular, when dealing with lengthy texts that
656 lack coherence, optimizing both token-level MI
657 and sequence-level MI maximization forces a se-
658 quence representation to resemble each individual
659 token embedding. The token-level MI maximiza-
660 tion objective provides no further improvement in
661 this case. This issue is exacerbated when some
662 terms are shared across clusters. This constraint
663 should be taken into account in future research.

664 Another limitation involving the general oper-
665 ation of contrastive learning is that the selection
666 of augmented samples directly affects the cluster-
667 ing performance. Notably, the best augmentation
668 strategy is still a subject of discussion and needs
669 more exploration. A study in Zhang et al. (2021)
670 and our own experiments with various augmenta-
671 tion settings show that varying an augmenter as
672 well as adjusting the configuration parameters both
673 affect the performance. Additionally, even if the
674 augmenter and the parameters used to generate
675 augmented texts are exactly the same, there is a
676 possibility that the outcomes from the two trials
677 may vary, adding a variance to the performance
678 results.

679 References

680 Sanjeev Arora, Yingyu Liang, and Tengyu Ma. 2017. A
681 simple but tough-to-beat baseline for sentence embed-
682 dings. In *5th International Conference on Learning
683 Representations, ICLR 2017, Toulon, France, April
684 24-26, 2017, Conference Track Proceedings*. Open-
685 Review.net.

686 Philip Bachman, R. Devon Hjelm, and William Buch-
687 walter. 2019. Learning representations by maximiz-
688 ing mutual information across views. In *Advances
689 in Neural Information Processing Systems 32: An-
690 nual Conference on Neural Information Processing
691 Systems 2019, NeurIPS 2019, December 8-14, 2019,
692 Vancouver, BC, Canada*, pages 15509–15519.

693 Mathilde Caron, Ishan Misra, Julien Mairal, Priya
694 Goyal, Piotr Bojanowski, and Armand Joulin. 2020.
695 Unsupervised learning of visual features by contrast-
696 ing cluster assignments. In *Advances in Neural In-*

*formation Processing Systems 33: Annual Confer-
697 ence on Neural Information Processing Systems 2020,
698 NeurIPS 2020, December 6-12, 2020, virtual*.
699

Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua,
700 Nicole Limtiaco, Rhomni St. John, Noah Con-
701 stant, Mario Guajardo-Cespedes, Steve Yuan, Chris
702 Tar, Yun-Hsuan Sung, Brian Strope, and Ray
703 Kurzweil. 2018. Universal sentence encoder. *CoRR*,
704 abs/1803.11175. 705

Ting Chen, Simon Kornblith, Mohammad Norouzi, and
706 Geoffrey E. Hinton. 2020. A simple framework for
707 contrastive learning of visual representations. In *Pro-
708 ceedings of the 37th International Conference on
709 Machine Learning, ICML 2020, 13-18 July 2020, Vir-
710 tual Event*, volume 119 of *Proceedings of Machine
711 Learning Research*, pages 1597–1607. PMLR. 712

Rotem Dror, Segev Shlomov, and Roi Reichart. 2019.
713 Deep dominance - how to properly compare deep
714 neural models. In *Proceedings of the 57th Annual
715 Meeting of the Association for Computational Lin-
716 guistics*, pages 2773–2785, Florence, Italy. Associa-
717 tion for Computational Linguistics. 718

John M. Giorgi, Osvald Nitski, Bo Wang, and Gary D.
719 Bader. 2021. Declutr: Deep contrastive learning for
720 unsupervised textual representations. In *Proceedings
721 of the 59th Annual Meeting of the Association for
722 Computational Linguistics and the 11th International
723 Joint Conference on Natural Language Processing,
724 ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual
725 Event, August 1-6, 2021*, pages 879–895. Association
726 for Computational Linguistics. 727

Jean-Bastien Grill, Florian Strub, Florent Alché,
728 Corentin Tallec, Pierre H. Richemond, Elena
729 Buchatskaya, Carl Doersch, Bernardo Ávila Pires,
730 Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal
731 Piot, Koray Kavukcuoglu, Rémi Munos, and Michal
732 Valko. 2020. Bootstrap your own latent - A new
733 approach to self-supervised learning. In *Advances
734 in Neural Information Processing Systems 33: An-
735 nual Conference on Neural Information Processing
736 Systems 2020, NeurIPS 2020, December 6-12, 2020,
737 virtual*. 738

Amir Hadifar, Lucas Sterckx, Thomas Demeester, and
739 Chris Develder. 2019. A self-training approach
740 for short text clustering. In *Proceedings of the
741 4th Workshop on Representation Learning for NLP
742 (RepLANLP-2019)*. Association for Computational
743 Linguistics. 744

Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and
745 Ross B. Girshick. 2020. Momentum contrast for un-
746 supervised visual representation learning. In *2020
747 IEEE/CVF Conference on Computer Vision and Pat-
748 tern Recognition, CVPR 2020, Seattle, WA, USA,
749 June 13-19, 2020*, pages 9726–9735. Computer Vi-
750 sion Foundation / IEEE. 751

R. Devon Hjelm, Alex Fedorov, Samuel Lavoie-
752 Marchildon, Karan Grewal, Philip Bachman, Adam
753

754	Trischler, and Yoshua Bengio. 2019. Learning deep representations by mutual information estimation and maximization . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	810
755		811
756		812
757		813
758		814
759	Sosuke Kobayashi. 2018. Contextual augmentation: Data augmentation by words with paradigmatic relations . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 2 (Short Papers)</i> , pages 452–457. Association for Computational Linguistics.	815
760		816
761		817
762		
763		818
764		819
765		820
766		821
767	Lingpeng Kong, Cyprien de Masson d’Autume, Lei Yu, Wang Ling, Zihang Dai, and Dani Yogatama. 2020. A mutual information maximization perspective of language representation learning . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	822
768		823
769		824
770		825
771		
772		826
773		827
774	Edward Ma. 2019. Nlp augmentation . https://github.com/makcedward/nlpaug .	828
775		829
776		830
777	Sosuke Nishikawa, Ryokan Ri, Ikuya Yamada, Yoshimasa Tsuruoka, and Isao Echizen. 2022. EASE: Entity-aware contrastive learning of sentence embedding . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3870–3885, Seattle, United States. Association for Computational Linguistics.	831
778		832
779		833
780		
781		834
782		835
783		836
784		837
785		838
786		839
787		840
788		
789		841
790		842
791		843
792		844
793		845
794		846
795		847
796		
797		848
798		849
799		850
800		851
801		852
802		853
803		854
804		855
805		
806		856
807		857
808		858
809		859
		860
		861
		862
		863
		864
		865
		866

867 *Lecture Notes in Computer Science*, pages 321–335.
868 Springer.

869 Jianhua Yin and Jianyong Wang. 2016. **A model-based**
870 **approach for text clustering with outlier detection**. In
871 *32nd IEEE International Conference on Data Engi-*
872 *neering, ICDE 2016, Helsinki, Finland, May 16-20,*
873 *2016*, pages 625–636. IEEE Computer Society.

874 Dessalew Yohannes and Yeregal Assabie. 2021.
875 **Amharic text clustering using encyclopedic knowl-**
876 **edge with neural word embedding**. *CoRR*,
877 abs/2105.00809.

878 Dejjiao Zhang, Feng Nan, Xiaokai Wei, Shang-Wen Li,
879 Henghui Zhu, Kathleen R. McKeown, Ramesh Nalla-
880 pati, Andrew O. Arnold, and Bing Xiang. 2021. **Sup-**
881 **porting clustering with contrastive learning**. *CoRR*,
882 abs/2103.12953.

883 Dejjiao Zhang, Wei Xiao, Henghui Zhu, Xiaofei Ma,
884 and Andrew Arnold. 2022. **Virtual augmentation**
885 **supported contrastive learning of sentence represen-**
886 **tations**. In *Findings of the Association for Com-*
887 *putational Linguistics: ACL 2022*, pages 864–876,
888 Dublin, Ireland. Association for Computational Lin-
889 guistics.

890 Xiang Zhang and Yann LeCun. 2015. **Text understand-**
891 **ing from scratch**. *CoRR*, abs/1502.01710.

892 Yan Zhang, Ruidan He, Zuozhu Liu, Kwan Hui Lim,
893 and Lidong Bing. 2020. **An unsupervised sentence**
894 **embedding method by mutual information maximiza-**
895 **tion**. In *Proceedings of the 2020 Conference on*
896 *Empirical Methods in Natural Language Process-*
897 *ing, EMNLP 2020, Online, November 16-20, 2020*,
898 pages 1601–1610. Association for Computational
899 Linguistics.

900 Xiaolin Zheng, Mengling Hu, Weiming Liu, Chaochao
901 Chen, and Xinting Liao. 2023. **Robust representation**
902 **learning with reliable pseudo-labels generation via**
903 **self-adaptive optimal transport for short text cluster-**
904 **ing**. In *Proceedings of the 61st Annual Meeting of the*
905 *Association for Computational Linguistics (Volume 1:*
906 *Long Papers)*, pages 10493–10507, Toronto, Canada.
907 Association for Computational Linguistics.

908 A Appendices

909 A.1 Datasets

910 Following previous works (Rakib et al., 2020;
911 Zhang et al., 2021; Pugachev and Burtsev, 2021;
912 Zheng et al., 2023), we conduct experiments and as-
913 sess the performance of our model on eight English
914 benchmark datasets for short text clustering. Table
915 2 presents the important statistics of all datasets.

- 916 • **AgNews**: a subset of the English news titles
917 dataset (Zhang and LeCun, 2015) in 4 differ-
918 ent topics, with 2,000 samples chosen ran-
919 domly from each topic by Rakib et al. (2020).

Dataset	$N^{Cluster}$	N^{Doc}	N^{Word}
AgNews	4	8,000	23
SearchSnippets	8	12,340	18
Biomedical	20	20,000	13
StackOverflow	20	20,000	8
Tweet	89	2,472	8
Googlenews-TS	152	11,109	28
Googlenews-T	152	11,109	6
Googlenews-S	152	11,109	22

Table 2: Dataset statistics. $N^{Cluster}$: number of clusters; N^{Doc} : number of short text documents; N^{Word} : average number of words in each document

- **SearchSnippets**: a dataset consisting of 12,340 web search snippets from 8 different categories (Phan et al., 2008).
- **Biomedical**: 20,000 paper titles, from 20 different Medical Subject Headings (MeSH), randomly selected by Xu et al. (2017) from the PubMed data distributed by BioASQ3.
- **StackOverflow**: challenge data published on Kaggle and randomly chosen by Xu et al. (2017), comprising 20,000 questions from Stack Overflow related to 20 distinct tags.
- **Tweet**: a dataset comprising 2,472 tweets with 89 groups (Yin and Wang, 2016).
- **GoogleNews**: GoogleNews-TS is a collection of titles and text snippets from 11,109 news articles covering 152 events (Yin and Wang, 2016). Only titles and snippet of each news article were extracted to produce GoogleNews-T and GoogleNews-S, respectively.

We spend up to 14 GPU hours on a Tesla V100 32G GPU to complete the training on all datasets for each MIST model’s configuration.

A.2 The Effects of Sequence- and Token-MI Maximization Objectives on NMI

Figure 4 shows the impact of the different ratios between the two MI maximization objectives on the clustering performance in terms of NMI across eight short text datasets. It follows the same trend as Accuracy as discussed in Section 5.2.1. MIST with our proposed generalized adaptive weighting function obtains the best clustering performance in terms of NMI for most datasets.

A.3 Positive Pairs in Contrastive Learning

It is a common practice in contrastive learning frameworks to only consider augmented texts as inputs, excluding original samples. However, we

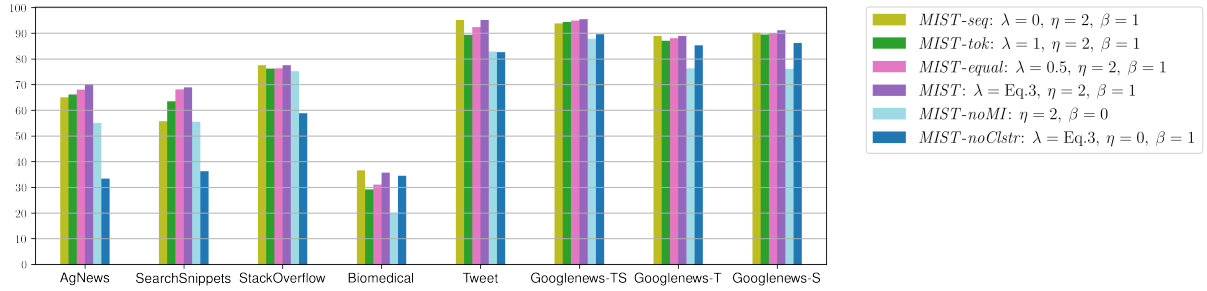


Figure 4: NMI for six different settings including four different weighting ratios between sequence-level and token-level MI maximization objectives. As well as, a setting where a clustering loss is absent ($\eta = 0$), and a setting where an MI loss is absent ($\beta = 0$). Note that when we set β to 0, λ has no effect.

	AgNews		SearchSnippets		StackOverflow		Biomedical	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
MIST w/ $\eta = 0$	56.96	33.40	50.30	36.30	64.40	58.80	43.26	34.55
MIST w/ $\eta = 1$	81.40	57.39	70.99	56.90	76.41	71.92	47.66	40.34
MIST w/ $\eta = 2$	89.47	70.25	76.72	67.69	78.74	77.59	39.15	34.66

	Tweet		GoogleNewsTS		GoogleNewsT		GoogleNewsS	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
MIST w/ $\eta = 0$	56.27	82.64	68.89	89.59	62.85	85.28	65.74	86.16
MIST w/ $\eta = 1$	64.46	86.27	74.86	91.89	66.91	87.04	71.98	88.58
MIST w/ $\eta = 2$	91.75	95.12	89.93	95.47	75.97	88.97	81.91	90.79

Table 3: The clustering results of MIST on three different weights of the clustering objective, η .

adopt a different input scheme. We discovered that feeding both original and augmented samples into our representation learning framework (as shown in Figure 1) yields better clustering results than exclusively taking two augmented texts as an input pair. One plausible reason is that when augmented texts are generated, the augmenter replaces some keywords in the original texts with new words. Short texts inherently have few keywords; hence, the absence of crucial words required for text categorization impacts clustering performance.

A.4 The Analysis of the Clustering Objective

As discussed in Section 5.2.2, the clustering performance is substantially affected by the weight of the clustering objective. Table 3 presents the performance of MIST across eight datasets in three situations, i.e., the coefficient of the clustering objective, η , in Eq.1 is assigned to 0, 1, and 2. The optimal results for the majority in terms of ACC and NMI are produced when η is set to 2.

A.5 Exploration of Data Augmentations

According to Zhang et al. (2021), which has studied the impacts of data augmentation in extensive details. The *Contextual Augmenter* has shown that it substantially outperforms other augmenters in their study. They hypothesized that since both the Con-

textual Augmenter and their encoder use the pre-trained transformers as the backbones, this allows the Contextual Augmenter to produce augmentation texts that are more informative and beneficial to their framework. We also adopted a pretrained transformer as the encoder in our framework and we observed that the experimental results followed the same trend as Zhang et al. (2021). We thus employ this augmenter in our experiments.

In this section, we investigate the impact of the Contextual Augmenter configurations in terms of masked language models and word substitution ratios. As shown in Table 4, we found that MIST using augmented texts generated from the BERT model with 20% substitution rate yields the best overall performance. Interestingly, MIST with augmented texts produced by *other masked language models with a 20% substitution rate* also yields outcomes close to those of BERT with the same substitution rate.

A.6 SCCL Reimplementation

To thoroughly compare the performance of our proposed representation learning strategy against the standard contrastive learning method in SCCL (Zhang et al., 2021), we reproduced SCCL in both an end-to-end version (SCCL) and a multiple-stage

	AgNews		SearchSnippets		StackOverflow		Biomedical	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
MIST w/ BERT 10%	87.74	66.99	75.98	67.71	77.78	76.42	37.51	33.97
MIST w/ BERT 20%	89.47	70.25	76.72	67.69	78.74	77.59	39.15	34.66
MIST w/ BERT 30%	86.33	66.09	81.46	67.71	73.60	71.55	39.79	34.61
MIST w/ RoBERTa 10%	87.51	66.81	75.64	67.11	77.84	76.50	38.61	35.11
MIST w/ RoBERTa 20%	88.85	69.12	76.21	68.52	77.74	76.41	37.17	31.62
MIST w/ RoBERTa 30%	86.43	66.4	73.77	65.72	77.76	77.03	29.48	27.38
MIST w/ DistilBERT 10%	87.22	66.44	74.96	65.89	77.67	76.30	38.29	34.29
MIST w/ DistilBERT 20%	89.42	70.26	75.74	67.85	77.72	77.05	38.29	32.31
MIST w/ DistilBERT 30%	87.96	67.66	74.23	64.11	77.67	76.34	38.83	34.63

	Tweet		GoogleNews-TS		GoogleNews-T		GoogleNews-S	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
MIST w/ BERT 10%	88.76	93.04	86.65	94.76	72.41	87.99	76.56	89.3
MIST w/ BERT 20%	91.75	95.12	89.93	95.47	75.97	88.97	81.91	90.79
MIST w/ BERT 30%	90.07	94.14	89.28	94.98	75.63	88.55	80.74	89.99
MIST w/ RoBERTa 10%	88.18	92.64	85.85	94.48	73.68	88.00	77.89	89.52
MIST w/ RoBERTa 20%	90.97	94.67	90.10	95.35	74.61	88.27	77.62	90.00
MIST w/ RoBERTa 30%	83.40	95.15	88.29	96.20	70.27	88.24	78.43	89.82
MIST w/ DistilBERT 10%	85.48	92.24	85.15	94.42	75.89	88.51	77.55	89.69
MIST w/ DistilBERT 20%	91.24	94.99	90.16	95.43	74.14	88.53	82.54	90.69
MIST w/ DistilBERT 30%	86.56	92.50	85.85	94.46	75.57	88.50	77.18	89.52

Table 4: The clustering performance of MIST when feeding augmented texts generated by Contextual Augmenter as inputs across nine different configurations.

version (SCCL-Multi). For the latter version, we apply the k -means algorithm on top of SCCL representations to make their pipeline identical to our framework except for the representation learning method. To be more specific, in this study, we report the experimental results of both reimplemented versions of SCCL using the *backbone* identified in the experimental setup of their publication. Moreover, SCCL considers the *Contextual Augmenter* with three configurations by setting the *word substitution ratio* of each text instance to 10%, 20%, and 30%. However, their study does not identify which setting produces the best outcomes. Therefore, we evaluate both reproduced versions of SCCL using *three alternative masked language models*: BERT-base, RoBERTa, and DistilBERT, with the aforementioned word substitution ratios for augmented pair generation to cover all scenarios reported in their study.

Table 5 reports the clustering performance of SCCL in both reproduced versions and in all configurations mentioned above. The reported performances show that despite the reproduced SCCL employing the configuration specified in their reference paper, their outcomes are still inferior to MIST in most cases. More specifically, MIST with the setup described in Section 4 outperforms SCCL and SCCL-Multi with the best parameter settings in the majority of cases. The fact that MIST produces better clustering performance than SCCL-Multi in

this study emphasizes that our proposed representation learning technique improves short text representations more effectively than the standard contrastive learning objective in the SCCL framework for short text clustering task. This demonstrates the success and efficiency of our proposed learning method even when compared with SCCL in various settings. Note that we collected the experimental results of reimplemented versions of SCCL from the *best iteration* for each dataset throughout 3000 iterations instead of using a stopping criterion, which is not indicated in their publication. Besides, the performances in their publication are reported from multiple settings.

Interestingly, the percentage of word replacement and masked language models employed for augmented text generation have an impact on the clustering performance. The best setting for these two parameters varies across different datasets. However, the performances of our proposed method presented in Table 1 are reported by using only a single setting for all datasets.

	AgNews		SearchSnippets		StackOverflow		Biomedical	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
SCCL (in the reference paper)	88.20	68.20	85.20	71.10	75.50	74.50	46.20	41.50
SCCL w/ BERT 10%	87.20	66.94	83.70	70.05	71.40	71.28	46.00	40.06
SCCL-Multi w/ BERT 10%	87.2	66.94	83.40	69.88	77.30	73.76	46.00	40.13
SCCL w/ BERT 20%	87.10	66.91	84.40	69.58	64.20	56.23	46.40	40.39
SCCL-Multi w/ BERT 20%	87.10	66.80	83.60	69.28	60.02	52.22	45.50	40.07
SCCL w/ BERT 30%	87.50	67.46	83.70	68.54	60.70	52.18	42.40	38.14
SCCL-Multi w/ BERT 30%	87.50	67.45	82.60	66.45	60.90	52.29	42.30	37.95
SCCL w/ RoBERTa 10%	87.00	66.57	84.50	70.21	62.10	54.26	28.50	20.35
SCCL-Multi w/ RoBERTa 10%	87.00	66.55	84.10	70.14	61.40	53.05	28.50	20.34
SCCL w/ RoBERTa 20%	85.20	64.20	62.60	41.66	60.70	52.26	39.60	32.66
SCCL-Multi w/ RoBERTa 20%	85.10	64.24	72.00	51.23	60.09	52.31	38.40	38.40
SCCL w/ RoBERTa 30%	84.00	62.24	30.70	10.07	60.70	52.28	39.10	32.77
SCCL-Multi w/ RoBERTa 30%	84.00	62.26	30.70	10.05	60.90	52.44	39.50	32.63
SCCL w/ DistilBERT 10%	87.30	67.16	84.70	70.79	70.20	69.49	46.10	39.87
SCCL-Multi w/ DistilBERT 10%	87.30	67.16	84.50	70.64	72.10	68.20	46.20	39.92
SCCL w/ DistilBERT 20%	86.80	65.87	84.70	70.62	71.40	69.38	46.30	39.94
SCCL-Multi w/ DistilBERT 20%	86.80	65.87	84.20	70.45	72.20	70.84	46.40	40.01
SCCL w/ DistilBERT 30%	87.20	66.77	85.00	71.63	70.80	70.04	46.30	40.49
SCCL-Multi w/ DistilBERT 30%	87.20	66.75	84.60	71.35	76.50	72.57	46.40	40.58

	Tweet		GoogleNews-TS		GoogleNews-T		GoogleNews-S	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
SCCL (in the reference paper)	78.20	89.20	89.80	94.90	75.80	88.30	83.10	90.40
SCCL w/ BERT 10%	56.80	81.91	70.10	89.49	62.50	81.53	69.00	86.29
SCCL-Multi w/ BERT 10%	75.30	88.39	86.70	93.95	76.30	88.25	81.00	89.82
SCCL w/ BERT 20%	57.10	82.54	75.60	90.99	63.00	81.72	67.80	85.97
SCCL-Multi w/ BERT 20%	78.20	89.41	88.70	94.70	76.20	87.97	81.10	89.60
SCCL w/ BERT 30%	56.6	82.23	74.2	90.83	61.30	81.20	64.9	89.78
SCCL-Multi w/ BERT 30%	78.80	89.58	89.90	94.91	75.60	87.88	82.10	89.77
SCCL w/ RoBERTa 10%	56.00	79.89	73.60	90.46	55.60	78.08	65.50	85.26
SCCL-Multi w/ RoBERTa 10%	71.10	85.86	86.60	93.94	56.90	78.52	80.50	89.50
SCCL w/ RoBERTa 20%	56.80	79.56	74.90	90.37	55.60	78.08	66.90	85.38
SCCL-Multi w/ RoBERTa 20%	74.20	86.61	88.10	94.27	58.40	79.28	81.30	89.87
SCCL w/ RoBERTa 30%	53.80	78.47	71.80	71.80	55.60	78.42	65.30	83.99
SCCL-Multi w/ RoBERTa 30%	63.60	76.98	85.20	93.53	56.60	78.42	78.00	88.14
SCCL w/ DistilBERT 10%	56.10	80.87	72.70	90.03	61.40	80.94	69.60	85.81
SCCL-Multi w/ DistilBERT 10%	78.80	88.91	87.70	94.25	74.30	87.78	79.70	89.20
SCCL w/ DistilBERT 20%	56.40	80.28	71.70	90.04	61.30	81.19	67.70	86.02
SCCL-Multi w/ DistilBERT 20%	77.10	88.61	86.50	94.03	75.10	87.51	79.50	89.70
SCCL w/ DistilBERT 30%	56.60	81.65	72.10	90.18	62.00	81.09	66.50	85.48
SCCL-Multi w/ DistilBERT 30%	76.00	88.39	88.50	94.18	75.80	87.60	79.10	89.01

Table 5: The clustering performances of the reimplemented SCCL and SCCL-Multi with nine different configurations for Contextual Augmenter. These configurations are obtained by setting the word substitution ratio of each text instance to 10% , 20%, and 30%, as well as using three alternative masked language models: BERT-base, RoBERTa, and DistilBERT.