

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 sample. Previous efforts based on dense representations are still inadequate since texts from different clusters are not sufficiently segregated in the embedding space prior to the clustering step. Even though the state-of-the-art technique integrated contrastive learning with a soft clustering objective to address this issue, the step in which all local tokens are summarized to form a sequence representation for the whole text may include noise that obscures the key information. We propose a framework called MIST: Mutual Information Maximization for Short Text Clustering, which overcomes the information limitation by maximizing the mutual information between text samples on both sequence and token levels. We assess the performance of our proposed method on eight standard short text datasets. Experimental results show that MIST outperforms the state-of-the-art methods 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. Hence, clustering these short texts is gaining more attention and becoming a crucial step for many real-world applications from recommendation to text retrieval (Yohannes and Assabie, 2021).

In short texts, words and phrases that are most representative 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 Bag-of-Words (BoW) and TF-IDF, provide relatively sparse representation vectors with limited descriptive power.

Hence, they perform poorly when clustered with a standard distance-based clustering method, such as k -means, in this situation (Hadifar et al., 2019).

To address this problem, deep neural networks have been employed to map high dimensional data into meaningful dense representations in a lower dimensional space. Most recent techniques for deep clustering follow a multi-phase style, in which the clustering process is carried out after learning feature representations (Xu et al., 2017; Hadifar et al., 2019; Yin et al., 2021). Unfortunately, the clustering performance of these methods remain unsatisfactory. One probable explanation is that texts still have a lot of overlap among categories in the latent space before clustering (Zhang et al., 2021).

Another deep clustering strategy optimizes representation learning and clustering objectives simultaneously (Zhang et al., 2021; Xie et al., 2016). To achieve desirable outcomes, Zhang et al. (2021) propose a method that adopts contrastive representation learning, which has been successful in self-supervised learning and is able to assist in spreading out the overlapped categories so that effective representations can be acquired, by simultaneously optimizing it along with a soft clustering target.

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 whole text representations while optimizing a contrasting objective. In particular, these representations are formed by summarizing all token representations in each text instance via mean pooling, including uninformative noises. We conjecture that this allows constructing a representation in which important information used to describe the text content may be obscured by noise, potentially affecting the clustering performance. Therefore, there is still a gap that needs to be explored in order to derive an efficient representation for short text clustering that does not omit informative terms.

084 In this paper, we introduce the **Mutual**
085 **Information Maximization Framework for Short**
086 **Text Clustering (MIST)**, a multi-stage framework
087 that learns textual representations by incorporating
088 two contrastive representation learning objectives
089 together with soft clustering assignments. Our con-
090 trastive learning procedure is based on mutual infor-
091 mation (MI) maximization, which facilitates us to
092 compare the semantic similarity across different hi-
093 erarchical levels to achieve multiple purposes. First,
094 we perform contrastive learning at a sequence-level
095 by contrasting between entire text representations.
096 Additionally, we also attempt to enforce each text
097 representation to extract information that is shared
098 across all of its tokens. In particular, we maximize
099 the MI between a text representation and all of its
100 local-level token embeddings to extract the shared
101 information among all the individual words in the
102 text. As a consequence, the information essential
103 to describe texts is preserved in the representations.

104 MIST handles the substantial challenge of short
105 text clustering, and our contributions are as follows:

- 106 • We propose MIST, a multi-stage framework
107 for short text clustering, which integrates two
108 contrastive learning objectives: (1) sequence-
109 level and (2) token-level MI maximization to
110 learn effective short text representations and
111 also be useful for clustering.
- 112 • To effectively balance sequence- and token-
113 level MI maximizations, we use a simple dy-
114 namic weighting function that adjust the ob-
115 jectives ratio in accordance with the length of
116 subword tokens in each minibatch.
- 117 • We conduct an extensive experiment to evalu-
118 ate the performance of MIST over eight
119 standard benchmarks of short text clustering.
120 MIST improves the clustering performance in
121 terms of Accuracy and NMI for most cases
122 compared to the current state-of-the-art.

123 2 Related Work

124 **Short Text Clustering.** There are a number of
125 approaches to overcome the sparsity of short text
126 representations, such as (1) multi-stage approaches
127 which break down the clustering framework into
128 multiple stages, (2) clustering enhancement algo-
129 rithms that apply outlier removal, and (3) a joint
130 framework that simultaneously optimizes both rep-
131 resentation learning and clustering objectives.

132 Several multi-stage works perform clustering
133 after learning feature representations. Pretrained-

134 word embeddings (Mikolov et al., 2013a,b; Pen-
135 nington et al., 2014) and neural networks are
136 adopted to transform data into meaningful repre-
137 sentations. Xu et al. (2015, 2017) use a convo-
138 lutional neural network to learn non-biased deep
139 feature representations by fitting the output units
140 with pretrained-binary codes from a dimensionality
141 reduction method. Hadifar et al. (2019) uti-
142 lize Smooth Inverse Frequency (SIF) (Arora et al.,
143 2017) to obtain weighted word embeddings. Dur-
144 ing training, they enrich discriminative features by
145 tuning an autoencoder with soft clustering assign-
146 ments from a clustering objective. For the afore-
147 mentioned works, k -means clustering is then em-
148 ployed on trained representations to get the final
149 clusters.

150 Another direction is to enhance the performance
151 of the initial clustering with an iterative classifica-
152 tion algorithm. Rakib et al. (2020) proposed an
153 ECIC algorithm which detects and removes out-
154 liers in each iteration. Moreover, they make use
155 of word embeddings by averaging them to repre-
156 sent each text, and combine the ECIC algorithm
157 with hierarchical clustering. To boost the cluster-
158 ing quality further, (Pugachev and Burtsev, 2021)
159 exploit deep sentence representations (Cer et al.,
160 2018) and made modifications to the ECIC algo-
161 rithm.

162 The recent state-of-the-art, SCCL (Zhang et al.,
163 2021), leverages a contrastive method from self-
164 supervised learning to encourage greater separa-
165 tion between overlapped categories in the original
166 data space. By jointly optimizing a contrastive loss
167 and a clustering objective (Reimers and Gurevych,
168 2019a), SCCL outperforms prior works and yields
169 cutting-edge results. In addition, other contrastive
170 learning methods have also been experimented on
171 short-text clustering, such as using entities for con-
172 trastive learning to provide supervision signals for
173 their related sentences (Nishikawa et al., 2022),
174 and using virtual augmentation for contrastive learn-
175 ing to circumvent the discrete nature of language
176 (Zhang et al., 2022). However, these methods do
177 not outperform SCCL on short text clustering.

178 **Self-supervised learning.** Self-supervision has
179 gained popularity and become a common technique
180 in unsupervised representation learning for a vari-
181 ety of downstream purposes. Many recent accom-
182 plishments have been based on contrastive repre-
183 sentation learning (Chen et al., 2020; He et al.,
184 2020; Caron et al., 2020; Grill et al., 2020).

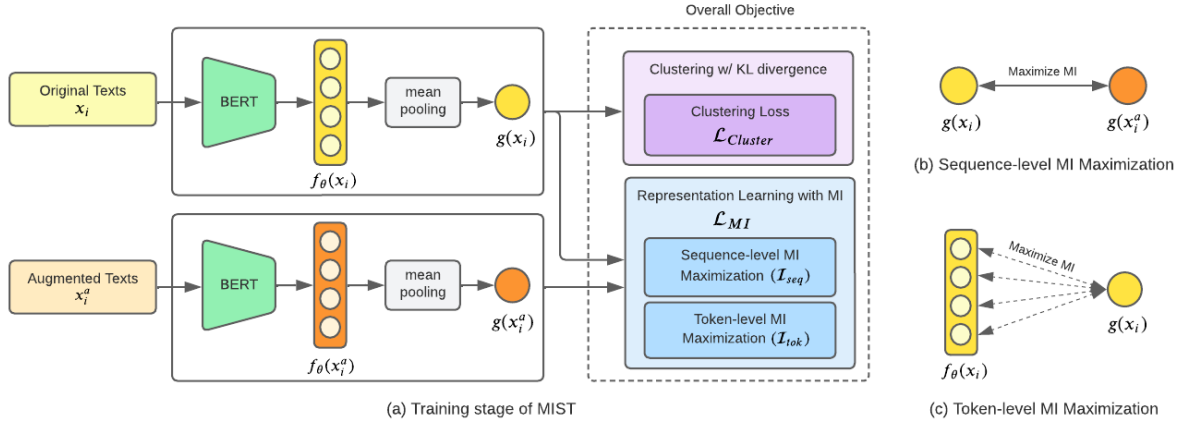


Figure 1: (a) The overview of the training stage of MIST. For each pair of original text x_i , and its augmented version x_i^a , MIST simultaneously optimizes the clustering objective $\mathcal{L}_{\text{Cluster}}$, and the MI maximization objectives \mathcal{L}_{MI} . The \mathcal{L}_{MI} comprises (b) a sequence-level MI maximization objective \mathcal{I}_{seq} , which attempts to maximize MI between sequence representations of x_i and x_i^a , along with (c) a token-level MI maximization objectives \mathcal{I}_{tok} that maximizes MI between a sequence representation (of both x_i and x_i^a) and its tokens ($f_\theta(x_i)$ and $f_\theta(x_i^a)$).

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. Each objective is then used separately according to the task. The authors also 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.

In this work, we leverage the MI maximization strategies to learn text representations specifically for short text clustering. We also investigate a weighting method for appropriately balancing MI objectives in order to improve clustering outcomes.

3 Proposed Method: MIST

In this section, we propose a short text clustering framework consisting of two steps: we first train a model using feature representation learning objectives as illustrated in Figure 1 and then apply the k -means clustering algorithm at inference time. The main idea of our solution lies in the proposed objective function \mathcal{L} that takes into account a MI

objective \mathcal{L}_{MI} that preserves a local invariance for each sample and an unsupervised clustering objective $\mathcal{L}_{\text{Cluster}}$ that captures categorical structure.

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

where β and η represent the trade-off between \mathcal{L}_{MI} , and $\mathcal{L}_{\text{Cluster}}$. We set β to 1, and η to 2 to give more weight to $\mathcal{L}_{\text{Cluster}}$.

We describe our proposed method in the following subsections. Section 3.1 provides a description for the MI maximization learning procedure, which includes (1) sequence-level and (2) token-level MI maximization objectives, along with a weighting function for balancing them. Section 3.2 presents the auxiliary clustering objective that enforces the encoder to create a suitable representation space for clustering.

3.1 Representation Learning with MI maximization

One strategy to improve clustering performance is to create an embedding space that minimizes local invariance for each individual sample via representation learning. A prominent method for creating such embedding space is contrastive learning which relies on contrasting representations throughout the whole context (global feature). Short text inputs are varied in terms of their lengths across different datasets. Consequently, there are short-text with smaller size (e.g., 6-8 words), as well as longer texts (e.g., 22-28 words). The latter tends to contain more words that may not be beneficial in defining high-level semantics useful for clustering.

Optimizing merely the global-level objective, as commonly done in contrastive learning, may not be sufficient to train effective representations for short text with weak signals problem.

To prevent local information from being obscured, we adopt an additional learning objective to constrain the representation of the entire text to contain high MI with each of its token embedding. In this investigation, we refer to the global and local features as *sequence* and *token* representations, respectively. Therefore, we build our training framework based on MI maximization strategy to reduce discrepancy between sequence- and token-level representations via their relative ability to predict each other across the representation levels.

Computing the MI Objective. As shown in Figure 1, the 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} .

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

where λ corresponds to the balancing weight for \mathcal{I}_{seq} and \mathcal{I}_{tok} objectives. We discovered that the number of tokens is an important factor in determining the ratio between the two objectives. In this study, we use a simple function to calculate the weight λ for each minibatch of size N depending on the length of each text:

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

and l_i denotes the number of tokens in a text x_i . Further discussion can be found in section 4.3.1.

In the learning stage, we first 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 includes a pretrained transformer network f_θ that 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. The sequence of token representations are then subsequently averaged by mean pooling operation $m(f_\theta(x))$ to generate a sequence representation denoted as $g(x) = m(f_\theta(x)) \in \mathbb{R}^d$.

Computing the Sequence-level MI. The first learning objective, \mathcal{I}_{seq} , aims to learn a representa-

tion that captures the entire context by contrasting samples at the sequence-level. According to Tian et al. (2020), contrastive learning is equivalent to maximizing the lower bound of MI between the representations of two texts. By treating each original text $g(x^o)$ and its augmentation $g(x^a)$ as positive samples, we can define \mathcal{I}_{seq} over the whole minibatch as follows.

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

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

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

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.

Computing the Token-level MI. To further enrich a text representation, we include a second learning objective, \mathcal{I}_{tok} , to MIST. Inspired by Zhang et al. (2020), this learning objective encourages a text representation to incorporate and preserve local information shared across all contextualized tokens. In particular, we attempt to maximize the average MI between a sequence representation and all of its token representations, while minimizing MI with the tokens of other texts. Conceptually, this reflects how much more precisely we can determine the representation of a token given a sequence representation compared to when we are unaware of the sequence representation (Bachman et al., 2019). We now define \mathcal{I}_{tok} for each minibatch as

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

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

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

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

3.2 Clustering with KL divergence

To encourage the coalescence of samples that are most likely to belong to the same cluster, we also employ a clustering objective $\mathcal{L}_{\text{Cluster}}$ along with the MI maximization objective. We follow the clustering method proposed by Xie et al. (2016), which are also used by Hadifar et al. (2019); Yin et al. (2021) and 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. In particular, 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 utilizing it to assist in refining clusters’ centroids. The main idea is to give more importance towards text samples with high clustering confidence. The probability $p_{jk} \in P$ is calculated as follows.

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

In order to match the soft cluster assignments to the target distribution, the KL-divergence between these two probability distributions, P and Q , is calculated 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 loss for each minibatch of size N as

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

4 Experiments

4.1 Experimental Setup

Datasets. We conduct experiments and evaluate the performance of MIST on the eight standard

short text clustering datasets, following previous works (Rakib et al., 2020; Zhang et al., 2021; Pugachev and Burtsev, 2021). Dataset descriptions and statistics are shown 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 the 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 select 10% of data as the validation set. Furthermore, we follow Zhang et al. (2021) by not performing any pre-processing operations on any of the eight datasets. Although some of existing works preprocessed the texts by removing symbols, stop words, punctuation or converting them to lowercase.

For the contrastive loss functions in the training stage, we consider original and augmented texts as inputs since we discovered that they are more effective than employing two augmented pairs in our experiment. To generate augmented samples for each text instance, we choose *Contextual Augmenter* (Kobayashi, 2018; Ma, 2019) using BERT and a 20% word substitution ratio. We found that this data augmentation setting can provide the best results as shown in Appendix A.6. We use two standard metrics, the clustering accuracy (ACC) and the normalized mutual information (NMI) to measure the clustering performance. The clustering accuracy is calculated via the Hungarian algorithm and the results are averaged over five trials.

4.2 Experimental Results

We compare the performance of our proposed framework, MIST, with state-of-the-art methods including STCC (Xu et al., 2017), Self-Train (Hadifar et al., 2019), HAC-SD (Rakib et al., 2020), SCA-AE (Yin et al., 2021) and SCCL (Zhang et al., 2021). As demonstrated in Table 1, MIST achieves state-of-the-art results for most cases in terms of Accuracy and NMI across eight benchmark datasets. In addition to the results reported in the reference papers, we further compare our method with SCCL, the state-of-the-art model that also employs contrastive learning for short text clustering, by reproducing SCCL in an end-to-end

	AgNews		SearchSnippets		StackOverflow		Biomedical	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
<i>Reported in the references</i>								
BoW [†]	27.6	2.6	24.3	9.3	18.5	14.0	14.3	9.2
TF-IDF [†]	34.5	11.9	31.5	19.2	58.4	58.7	28.3	23.2
Skip-Thought [‡]	-	-	33.6	13.8	9.3	2.7	16.3	10.7
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
SCCL [†]	88.2	68.2	85.2	71.1	75.5	74.5	46.2	41.5
<i>Reimplement</i>								
SCCL w/ BERT 20%	87.10	67.18	84.78	70.02	49.48	47.50	44.90	39.73
SCCL-Multi w/ BERT 20%	86.95	67.06	83.88	69.50	53.56	46.99	44.70	39.65
<i>Proposed Method</i>								
MIST	89.47	70.25	76.72	67.69	78.74	77.59	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>								
BoW [†]	49.7	73.6	57.5	81.9	49.8	73.2	49.0	73.5
TF-IDF [†]	57.0	80.7	68.0	88.9	58.9	79.3	61.9	83.0
Skip-Thought [‡]	-	-	-	-	-	-	-	-
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
SCCL [†]	78.2	89.2	89.8	94.9	75.8	88.3	83.1	90.4
<i>Reimplement</i>								
SCCL w/ BERT 20%	55.98	82.12	75.35	90.96	62.53	81.95	67.88	86.07
SCCL-Multi w/ BERT 20%	79.05	89.59	88.83	94.69	76.20	87.89	82.25	90.01
<i>Proposed Method</i>								
MIST	91.75	95.12	89.93	95.47	75.97	88.97	81.91	90.79

Table 1: Experimental results on eight short text clustering datasets. [†] and [‡] refer to results taken from Zhang et al. (2021) and Hadifar et al. (2019), respectively; both originally present their results in one decimal place.

(original) as well as a multi-stage version analogous to our architecture for fair comparison. The reimplemented versions of SCCL employ the same augmentation setting as our model. We refer to these model as SCCL w/ BERT 20% and SCCL-Multi w/ BERT 20%, respectively. The comparative results in Table 1 show that MIST outperforms SCCL, SCCL w/ BERT 20% and SCCL-Multi w/ BERT 20% in 11, 12 and 10 cases, respectively.

For datasets with small number of clusters, Search Snippets and Biomedical, MIST does not yield competitive results. We obtain a weaker result on Biomedical, since the dataset used to pre-train our encoder is a general domain one. On the other hand, Hadifar et al. (2019) produces the best result using pretrained embeddings learned from a large in-domain biomedical corpus. For the SearchSnippets dataset, MIST also obtains a poorer result. One probable explanation is that snippets are typically composed of content words, as well as the dataset has been automatically crawled and preprocessed further by Phan et al. (2008), the pre-

processing steps include removing stop and rare words. Due to the length and incoherency of each text in this dataset, our algorithm becomes dependent on keywords rather than contextual information. Particularly, when it performs the token-level MI maximization objective in the representation learning stage, which enforces similarity between each contextualized token representation and the sequence representation of the incoherent text sequence. This can be problematic when the same keywords also appear in different clusters.

For datasets with a large number of clusters, such as GoogleNews, it is more likely that texts in different clusters may share a similar content due to fine-grained categorization, inducing ambiguity. We conjecture that this ambiguity in textual data and ground-truths is causing inaccurate predictions. As GoogleNews-T only contains news headlines, which are relatively short with few keywords. It presents a challenge for clustering these texts into a large number of categories. For example, "liam adam sentenced abuse daughter" is a

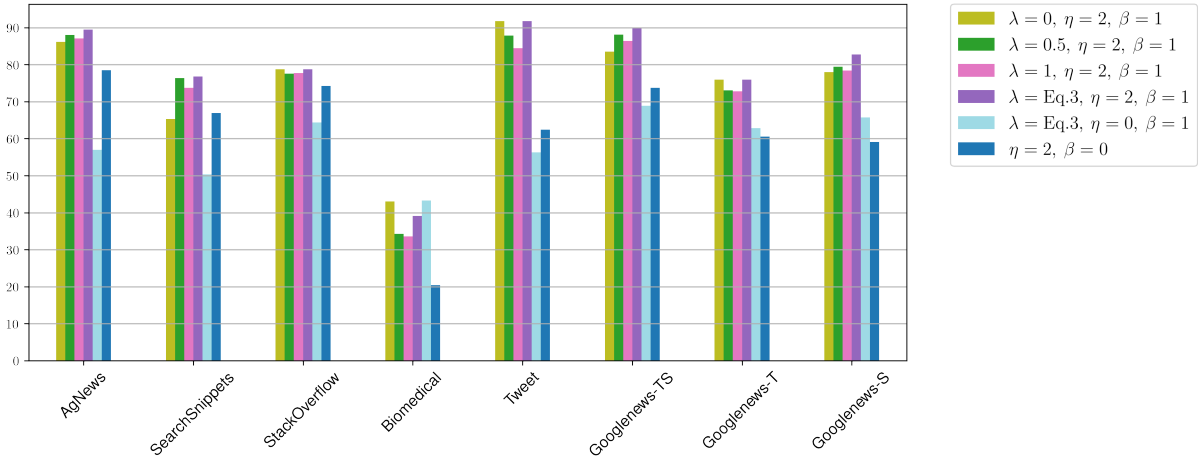


Figure 2: Accuracy for six different settings including four different weighting ratios between sequence- and token-level MI maximization objectives. As well as, a setting where a clustering loss is absent ($\eta = 0$), and a setting where an MI loss is absent ($\beta = 0$). Note that when we set β to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

news headline in a cluster of news related to Gerry Adams, an IRA activist and the former president of Sinn Féin. This sample contains same keywords found in another cluster with news about domestic violence. Another cause of inaccuracy is when the content of texts in one cluster is a subtopic of the content in another cluster.

We hypothesize that Rakib et al. (2020), which employs hierarchical clustering and outlier removal algorithms, can better deal with hierarchical nature of data. Consequently, Rakib et al. (2020) outperforms our method and SCCL on this scenario in terms of Accuracy on this dataset. While our method and SCCL both aim to improve representations through the use of contrastive representation learning. As shown in Table 1, MIST also has lower Accuracy on GoogleNews-T and GoogleNews-S than the reported result of SCCL in the reference paper and SCCL-Multi w/ BERT 20%, respectively. Where we collected the experimental results of SCCL w/ BERT 20% and SCCL-Multi w/ BERT 20% from the best iteration for each dataset instead of using a stopping criterion, which is also not mentioned in Zhang et al. (2021).

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. MIST can derive a very strong and comparable Accuracy to SCCL on GoogleNews-S and outperforms SCCL on GoogleNews-TS. This is because, GoogleNews-S contains text snippets of Google News, and

GoogleNews-TS includes both the titles and snippets.

Additional details and the comparison results of SCCL in both reproduced versions with other augmentation settings can be found in the A.5. According to the results in A.5, our method still outperforms SCCL in both end-to-end and multi-phase settings in 11 cases.

4.3 Ablation Study

To better understand the impact of each component in our training procedure on the clustering performance, we conduct additional experiments by varying the ratio setting between sequence- and token-level MI maximization objectives in the MI loss \mathcal{L}_{MI} , as well as the clustering objective $\mathcal{L}_{Cluster}$.

4.3.1 The effects of sequence- and token-MI maximization objectives

Let us consider the effects of sequence- and token-level MI maximization objectives on the clustering performance. We report the performance of our model in four different ratios by setting λ in Eq.2 to 1, 0.5, 0, and also assigning the value to λ using Eq. 3. In this section, we refer to the MIST model with a sequence-only MI maximization ($\lambda = 0$) and a token-only ($\lambda = 1$) MI maximization objectives as MIST-seq and MIST-tok, respectively. As demonstrated in Figure 2, MIST with the ratio set according to Eq.3 yields the best performance in terms of Accuracy, except for Biomedical. NMI tends to follow the same direction as Accuracy, as demonstrated in Appendix A.2. This indicates that

the length of texts (the amount of token embeddings) is a major consideration in the selection of appropriate ratios between both MI maximization objectives. In addition, we also investigate the scenario when both MI objectives are absent ($\beta = 0$). The ablation results reveal that when both MI maximization objectives are removed, the performance suffers substantially on all datasets. This shows that the MI loss is necessary for performance gain.

For datasets with long-length texts, such as GoogleNews-TS, we discovered that MIST produces the best outcomes when token- and sequence-level MI maximization objectives are weighted using λ calculated by Eq. 3. Note that this setting also outperforms the scenario when both objectives are assigned the same weight ($\lambda = 0.5$). We can also see that MIST-tok always outperforms MIST-seq. This shows that if the text is lengthy, MIST-seq may not be sufficient. This is because informative terms of the text are averaged with other non-informative terms via mean pooling. Since infrequent keywords in the text are more likely to be overlooked, maximizing each local token embeddings with its sequence representation helps alleviate this problem.

For datasets with very short-length texts, such as StackOverflow and Tweet, the weighting ratio based on Eq. 3 is equivalent to setting λ to 0. In this situation, MIST is identical to MIST-seq. MIST-seq outperforms other settings, followed by MIST with integrating the seq- and token-level MI maximization objectives which always performed better than MIST-tok. For instance, texts in the Tweet dataset are relatively short and contains solely content words rather than coherent texts. As a result, our model with token-level MI maximization objective, MIST-tok and MIST with the combination of token- and sequence-MI maximization objectives, might emphasize on keywords that could also appear in multiple clusters, causing ambiguity.

4.3.2 The effects of soft cluster assignments

As shown in Figure 2, the clustering performance drops significantly when we remove the clustering with KL divergence objective ($\eta = 0$). This demonstrates that the categorical structure imposed by simultaneously optimizing the clustering loss with the representation learning objectives is a crucial component that boosts performance. However, this trend holds true for all datasets, except for Biomedical. One possible explanation is that, since the encoder was not pretrained with textual infor-

mation which was suitable for its specific domain, the clustering loss does not benefit the efficiency of our model than the representation objectives.

Furthermore, we observe that as the weight of clustering increases, the performance continuously improves until it reaches saturation as η , the weight for the clustering loss, approaches 2. As depicted in Figure 3, the accuracy and NMI of AgNews both improve as we gradually increase the clustering weight until the appropriate value, which is 2 in our experiment. The supplementary experimental results can be found in Appendix A.4.

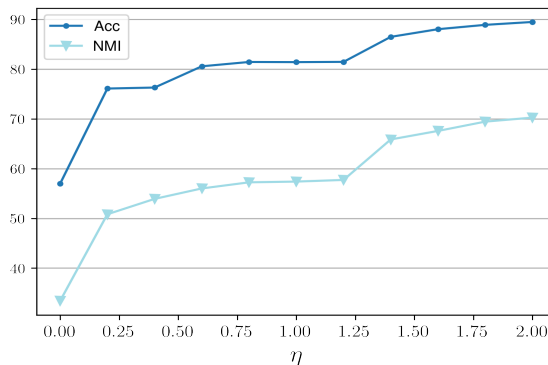


Figure 3: The clustering performance on AgNews based on the strength of the clustering loss. The strength of both MI maximization objectives are kept constant based on Eq. 3

5 Conclusion

We propose a novel multi-stage framework that employs two contrastive learning objectives based on MI maximization methods to produce effective representations for short texts. To learn distinct text representations, the first contrastive learning objective maximizes MI between original texts and their augmentations at the sequence level. And the second objective maximizes MI between sequence representations and their local tokens. Additionally, we introduce a preliminary weighting function for properly balancing the two MI maximization objectives during training process.

We have conducted extensive experiments across eight benchmark datasets for short text to study the effectiveness of our method. Our model outperforms state-of-the-art methods in most cases on Accuracy and NMI. This demonstrates that utilizing the MI maximization strategy during the contrastive learning process could potentially be a promising tactic. Further study would be worthwhile since it might enhance the quality of textual representations for other tasks

References

- 617 Sanjeev Arora, Yingyu Liang, and Tengyu Ma. 2017. [A simple but tough-to-beat baseline for sentence embeddings](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net. 674
- 618
- 619
- 620
- 621
- 622
- 623 Philip Bachman, R. Devon Hjelm, and William Buchwalter. 2019. [Learning representations by maximizing mutual information across views](#). In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 15509–15519. 676
- 624
- 625
- 626
- 627
- 628
- 629
- 630 Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. 2020. [Unsupervised learning of visual features by contrasting cluster assignments](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. 677
- 631
- 632
- 633
- 634
- 635
- 636
- 637 Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. [Universal sentence encoder](#). *CoRR*, abs/1803.11175. 678
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647
- 648
- 649
- 650
- 651
- 652
- 653
- 654
- 655
- 656
- 657
- 658
- 659
- 660
- 661
- 662
- 663
- 664
- 665
- 666
- 667
- 668
- 669
- 670
- 671
- 672
- 673
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 722
- 723
- 724
- 725
- 726
- 727
- 728
- 4th Workshop on Representation Learning for NLP (RepL4NLP-2019). Association for Computational Linguistics.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. 2020. [Momentum contrast for unsupervised visual representation learning](#). In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 9726–9735. Computer Vision Foundation / IEEE.
- R. Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Philip Bachman, Adam Trischler, and Yoshua Bengio. 2019. [Learning deep representations by mutual information estimation and maximization](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Sosuke Kobayashi. 2018. [Contextual augmentation: Data augmentation by words with paradigmatic relations](#). In *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)*, pages 452–457. Association for Computational Linguistics.
- 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 *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Edward Ma. 2019. Nlp augmentation. <https://github.com/makcedward/nlpaug>.
- Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. [Efficient estimation of word representations in vector space](#). In *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*.
- Tomás Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013b. [Distributed representations of words and phrases and their compositionality](#). In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pages 3111–3119.
- Sosuke Nishikawa, Ryokan Ri, Ikuya Yamada, Yoshimasa Tsuruoka, and Isao Echizen. 2022. [EASE: Entity-aware contrastive learning of sentence embedding](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3870–3885, Seattle, United States. Association for Computational Linguistics.

729	Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. 2016. f-gan: Training generative neural samplers using variational divergence minimization . In <i>Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain</i> , pages 271–279.	786
730		787
731		788
732		789
733		790
734		791
735		792
736	Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch . In <i>NIPS 2017 Workshop on Autodiff</i> .	793
737		794
738		795
739		796
740		797
741	Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL</i> , pages 1532–1543. ACL.	798
742		799
743		
744		800
745		801
746		802
747		803
748	Xuan Hieu Phan, Minh Le Nguyen, and Susumu Horiguchi. 2008. Learning to classify short and sparse text & web with hidden topics from large-scale data collections . In <i>Proceedings of the 17th International Conference on World Wide Web, WWW 2008, Beijing, China, April 21-25, 2008</i> , pages 91–100. ACM.	804
749		805
750		806
751		807
752		
753		808
754		809
755	Leonid Pugachev and Mikhail S. Burtsev. 2021. Short text clustering with transformers . <i>CoRR</i> , abs/2102.00541.	810
756		811
757		
758	Md. Rashadul Hasan Rakib, Norbert Zeh, Magdalena Jankowska, and Evangelos E. Milios. 2020. Enhancement of short text clustering by iterative classification . In <i>Natural Language Processing and Information Systems - 25th International Conference on Applications of Natural Language to Information Systems, NLDB 2020, Saarbrücken, Germany, June 24-26, 2020, Proceedings</i> , volume 12089 of <i>Lecture Notes in Computer Science</i> , pages 105–117. Springer.	812
759		813
760		814
761		815
762		816
763		817
764		818
765		819
766		820
767	Nils Reimers and Iryna Gurevych. 2019a. Sentencebert: Sentence embeddings using siamese bert-networks . <i>CoRR</i> , abs/1908.10084.	821
768		822
769		823
770	Nils Reimers and Iryna Gurevych. 2019b. Sentencebert: Sentence embeddings using siamese bert-networks . 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 2019, Hong Kong, China, November 3-7, 2019</i> , pages 3980–3990. Association for Computational Linguistics.	824
771		825
772		
773		826
774		827
775		828
776		829
777		
778	Alessandro Sordani, Nouha Dziri, Hannes Schulz, Geoffrey J. Gordon, Philip Bachman, and Remi Tachet des Combes. 2021. Decomposed mutual information estimation for contrastive representation learning . In <i>Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event</i> , volume 139 of <i>Proceedings of Machine Learning Research</i> , pages 9859–9869. PMLR.	830
779		831
780		832
781		833
782		834
783		
784		835
785		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

Xiang Zhang and Yann LeCun. 2015. [Text understanding from scratch](#). *CoRR*, abs/1502.01710.

Yan Zhang, Ruidan He, ZuoZhu Liu, Kwan Hui Lim, and Lidong Bing. 2020. [An unsupervised sentence embedding method by mutual information maximization](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 1601–1610. Association for Computational Linguistics.

A Appendices

A.1 Datasets

Following previous works, we conduct experiments and evaluate the performance of our model on the eight short text clustering datasets. These datasets only contain texts in English. All of them are publicly available online. A summary of the statistics of all datasets is listed in Table 2.

- **AgNews**: a subset of the dataset of English news titles (Zhang and LeCun, 2015) across 4 different topics, where 2,000 samples from each topic were randomly chosen by Rakib et al. (2020).
- **SearchSnippets**: a dataset comprising 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), which consists of 20,000 question titles from Stack Overflow related to 20 distinct tags.
- **Tweet**: a dataset comprising 2,472 tweets with 89 groups (Yin and Wang, 2016).
- **GoogleNews**: a collection of both titles and text snippets from 11,109 news articles covering 152 events (Yin and Wang, 2016). Only the titles and the text snippet of each news article were extracted out of the GoogleNews-TS 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.

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

A.2 The effects of sequence- and token-MI maximization objectives on NMI

Figure 4 shows the effects of sequence- and token-MI maximization objectives on NMI. It follows the same trend as Accuracy as discussed in Section 4.3.1.

A.3 Positive Pairs in Contrastive Learning

It is a common practice in contrastive learning frameworks to only consider augmented data as inputs, excluding an original sample. However, we adopt a different input scheme. We discovered that feeding both original and augmented samples into our contrastive learning framework (as shown in Figure 1) yields better clustering results than exclusively taking two augmented texts as an input pair. One probable explanation is that when augmented texts are created, the augmentser replaces some keywords in original texts with new words. Since short texts are typically short and include few keywords, the absence of crucial words required for text categorization has an impact on clustering performance.

A.4 The impact of soft cluster assignments

As discussed in Section 4.3.2, the clustering performance is substantially affected by varying the weight of the clustering objective during training representations process. 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 provided by MIST when η is set to 2.

A.5 SCCL Reimplementation

To thoroughly compare the performance of our contrastive learning strategy against SCCL, an existing

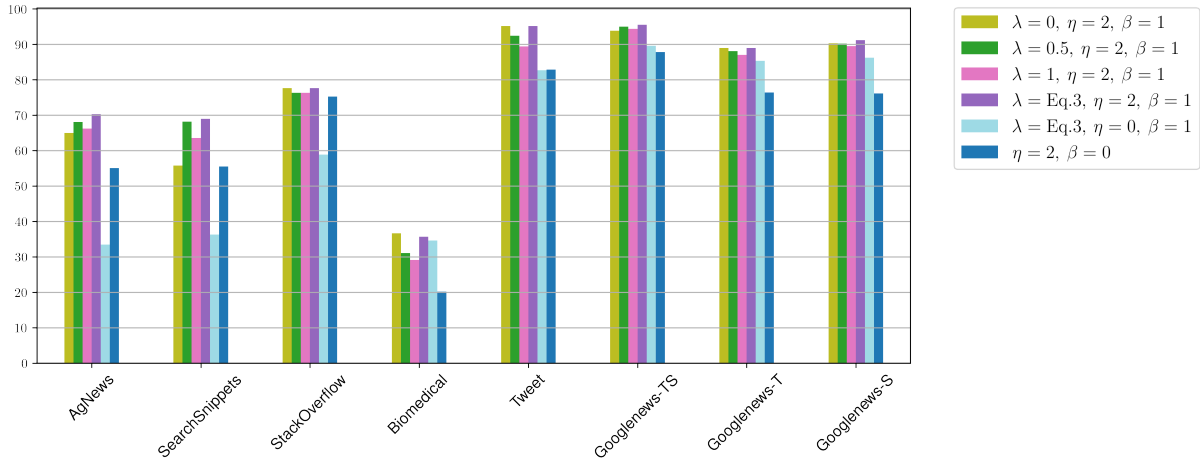


Figure 4: NMI for six different settings including four different weighting ratios between sequence- and token-level MI maximization objectives. As well as, a setting where a clustering loss is absent ($\eta = 0$), and a setting where an MI loss is absent ($\beta = 0$). Note that when we set β to 0, 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 .

924 contrastive learning method for short-text clustering, we reproduced SCCL in both original version
 925 and a multiple-stage version (SCCL-Multi), by applying the k -means algorithm on top of SCCL rep-
 926 resentations to make their pipeline identical to our framework. We followed Zhang et al. (2021) and
 927 used *Contextual Augmenter*, which was reported to offer the best results, to generate augmented texts
 928 for all training frameworks in this experiment. In the reference paper, SCCL considers *Contextual*
 929 *Augmenter* with three configurations by setting the word substitution ratio of each text instance to 10%,
 930 20%, and 30%. But their study does not identify which configuration for *Contextual Augmenter* set-
 931 ting produces the best outcomes. Therefore, we examine SCCL-Multi with three alternative masked
 932 language models: BERT-base, RoBERTa and DistilBERT for augmented pairs generation to covers
 933 all scenarios.

943 Table 4 reports the best clustering results for SCCL and SCCL-Multi in all configurations ob-
 944 tained during maximum iteration, as well as the

946 best results for SCCL produced using the Contextual Augmenter presented in Zhang et al. (2021).
 947 The percentage of word replacement and masked language models employed for augmented text gen-
 948 eration have an impact on the clustering performance of SCCL-Multi, since the best setting for
 949 these two parameters varies across datasets. Our contrastive learning approach outperforms both
 950 SCCL-Multi and SCCL with the best augmentation parameters settings in 6 out of 8 datasets.
 951

956 A.6 Exploration of Data Augmentations

957 According to Zhang et al. (2021), we investigate the impact of the *Contextual Augmenter* configu-
 958 rations in terms of masked language models and substitution percentage, respectively. As shown
 959 in Table 5, we found that MIST using augmented texts generated from the BERT model with 20%
 960 substitution rate during training step yields the best overall performance. MIST with augmented texts
 961 produced by other encoders with 20% substitution rate also yield the outcomes close to those of BERT
 962

967 with the same substitution rate.

968 **A.7 Limitations**

969 Despite the state-of-the-art performance, there are
970 several limitations, which we highlight in this sec-
971 tion. Firstly, the backbone of our model is pre-
972 trained using general domain data. Hence, when
973 our model encounters short texts in a specific do-
974 main, such as Biomedical, the performance drops
975 drastically. Furthermore, our representation learn-
976 ing procedure also performs poorly on short texts
977 with only content words or incoherent text se-
978 quences. Learning representations for incoherent
979 texts, by incorporating token-level MI maximiza-
980 tion objective, forces a sequence representation to
981 resemble each individual token embedding. For
982 short-texts with incoherent text, the token-level MI
983 maximization objective gives no further improve-
984 ment. This constraint should be taken into account
985 in future research.

986 Another limitation of our framework is that aug-
987 mented samples are crucial for the learning process
988 according to the general operation principle of con-
989 trastive learning. However, the best augmentation
990 strategy is still a subject of discussion and explo-
991 ration. A study in SCCL and comparison results
992 of our model with several augmentation settings
993 demonstrate that varied augments as well as dif-
994 ferent configuration factors have an on clustering
995 performance. Additionally, even if the technique
996 and the parameters used to generate augmented
997 texts are exactly the same, there is a possibility that
998 the outcomes from the two trials may vary, adding
999 a variance to the performance results.

	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 4: 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.

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