

# IOPC: Aligning Semantic and Cluster Centers for Few-shot Short Text Clustering

Anonymous ACL submission

## Abstract

In short text clustering, determining appropriate cluster centers is crucial. However, due to the limitations of short text representation quality, traditional methods often struggle to align cluster center with the core concept (*semantic center*) of each category, degrading clustering performance. To address this issue, we propose **IOPC**, a novel few-shot learning framework that achieves the alignment through two key modules. First, to capture effective semantics, we introduce an Interaction-enhanced Optimal Transport (**IEOT**) that leverages semantic interactions between samples to generate confident pseudo-labels. Based on these high-quality pseudo-labels, pseudo-semantic centers (*prototypes*) can be obtained. Furthermore, we propose Prototype-based Contrastive Learning (**PBCL**) to optimize text representations towards their corresponding prototypes. As training progresses, the continuous updating of pseudo-labels and prototypes gradually reduces the gap between cluster centers and semantic centers, improving clustering performance and stability. Extensive experiments on eight benchmark datasets show that IOPC outperforms state-of-the-art methods, achieving up to 7.34% improvement in accuracy on challenging Biomedical datasets and excelling in clustering stability and efficiency. The code is available at: <https://anonymous.4open.science/r/IOPC-origin-2522>.

## 1 Introduction

Short text clustering, which groups short texts into distinct clusters based on their semantic similarity, has broad applications in real-world domains such as chatbots (Kuhail et al., 2023), topic discovery (Murshed et al., 2023), and spam detection (Abkenar et al., 2023; Teja Nallamothu and Shais Khan, 2023). A key factor in achieving high-quality clustering is determining the appropriate cluster center for each category, as this critically influences whether samples can be grouped according to their

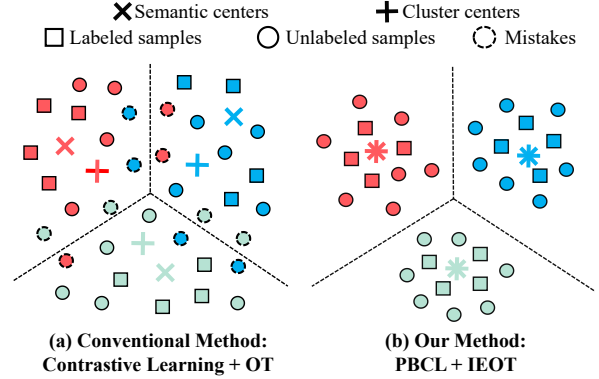


Figure 1: **Schematic Illustration of the Motivation.** (a) Previous works generate cluster centers that are misaligned with the underlying semantic centers. (b) In contrast, our method leverages few-shot samples effectively align cluster centers with the semantic centers, thereby facilitating accurate category aggregation.

intrinsic semantic similarities (Bai et al., 2012). The ideal scenario is that the cluster center for each category precisely corresponds to the core concept (*semantic center*) of that category in feature space. However, due to the lack of labeled samples and limitations in text representation quality, extracting the semantic center of each category remains a challenge (Fini et al., 2023). As illustrated in Figure 1(a), cluster centers often fail to align with the true semantic centers, leading to suboptimal category aggregation.

Previously, Zheng et al. (2023); Li et al. (2024) proposed constructing pseudo-labels to assign preliminary category information to certain samples, allowing similar samples to gradually converge during the iterative process. However, ensuring the quality of pseudo-labels remains a challenge. On the other hand, to enhance the quality of the text representation, Zhang et al. (2021); Chen et al. (2020) introduced contrastive learning, which optimizes text representations by pulling positive samples together and pushing negative samples apart in the feature space. However, these methods only

consider instance-level representation relationships, causing samples that should belong to the same category to be pushed apart, affecting cluster quality (Wang and Isola, 2020).

In this work, we propose **IOPC**, a novel few-shot short text clustering framework. In the optimization process of the feature space, it imposes constraints to pulling the text representation toward their semantic center. IOPC combines two key components: Interaction-enhanced Optimal Transport (**IEOT**) and Prototype-based Contrastive Learning (**PBCL**).

Specifically, (1) we integrate similarity interactions between samples into the optimal transport (OT) framework, enabling IEOT to generate more reliable pseudo-labels. We then seamlessly combine minimal labels with these pseudo-labels to effectively design a pseudo-semantic center (*prototype*) for each category. (2) Subsequently, we leverage these prototypes as targets, guiding samples of the same category to move closer to their corresponding prototype while distancing themselves from other prototypes through PBCL. As training progresses, the continuous updating of pseudo-labels drives the prototypes to gradually approach the true semantic centers, which in turn guides the text representations to move closer to them. Ultimately, IOPC will achieve alignment between the semantics of the cluster with the semantic center, as illustrated in Figure 1(b).

We demonstrate that IOPC achieves state-of-the-art performance on eight benchmark datasets. Notably, IOPC achieved the highest accuracy in all datasets, with improvements exceeding **7.34%** and **4.18%** on Biomedical and GoogleNews-T, respectively. Additionally, we show that our method exhibits faster convergence and more robust training compared to current methods. In summary, our main contributions are as follow:

- We propose a few-shot framework, IOPC, which integrates the following two key components, bridging the gap between semantic and cluster centers.
- We propose a novel optimal transport strategy, IEOT, which leverages interactions between samples. It generates reliable pseudo-labels to help the few-shot labels uncover the true semantic centers of each category.
- We propose a novel contrastive learning method, PBCL, which aligns cluster centers

with semantic centers by leveraging prototypes of core concepts to guide representation optimization.

- IOPC shows state-of-the-art results on eight benchmark datasets. Furthermore, our method achieves faster convergence and better stability compared to previous methods.

## 2 Related Works

**Short Text Clustering.** Short text clustering is challenging due to the limited number of words in short texts. In recent years, deep joint clustering methods have become mainstream by integrating representation learning and clustering into a unified framework. Notable examples include SCCL (Zhang et al., 2021), which uses DEC (Xu et al., 2017) as the clustering objective and contrastive learning to guide representation learning. RSTC (Zheng et al., 2023) proposes the use of pseudo-labels to assist the model in learning sample representations and clustering. STSPL-SSC (Nie et al., 2024) is built on the RSTC method, using fewer labeled data to assist the pseudo-labeling process. COTC (Li et al., 2024) combines sentence-level and token-level information to achieve more efficient clustering.

**Few-shot learning.** Few-shot methods leverage a small amount of labeled data and a large collection of unlabeled data to train models. The most intuitive approach is Pseudo-Labels (Lee et al., 2013), where a model trained on labeled data generates pseudo-labels for unlabeled examples, which are then added to the labeled set for the next iteration. However, hard labels easily exacerbate the classification bias of the training model (confirmation bias) (Arazo et al., 2020). To counteract this issue, researchers have shown benefits from soft labels and confidence thresholding (Arazo et al., 2020) as well as from different training strategies like co- and tri-training (Dong-DongChen and WeiGao, 2018; Nassar et al., 2021). In our research, we integrate optimal transport and pseudo-labeling methods to explore textual features and similarities, maximizing the guiding role of labeled information.

**Contrastive Learning.** As a promising paradigm of unsupervised learning, contrastive learning has lately achieved state-of-the-art performance in many fields (Grill et al., 2020). Contrastive learning aims to map data to a feature space where positive pairs are similar and negative pairs are dissimilar (Hadsell et al., 2006). Recently, (Zhang et al.,

2021) applies contrastive learning to short text clustering, upon which methods like (Zheng et al., 2023; Nie et al., 2024; Li et al., 2024) and many others have introduced further improvements. The previous methods typically distribute the samples uniformly in feature space (Wang and Isola, 2020), whereas our approach further optimizes them by incorporating semantics, thereby achieving consistency and accuracy.

### 3 Method

In this paper, we propose **IOPC** which is primarily attributed to two key factors: Interaction-enhanced Optimal Transport (**IEOT**) and Prototype-based Contrastive Learning (**PBCL**). An overall illustration of IOPC is shown in Figure 2. Specifically, the probability distributions obtained after the samples pass through the Encoder and Classifier are processed by IEOT to generate pseudo-labels. Subsequently, Prototypes are constructed from labeled samples and confident pseudo-labeled samples. PBCL makes the projections of the samples in the feature space to be more compact towards their corresponding prototypes, ensuring stable alignment between cluster centers and the true semantic centers. Subsequently, Instance-level Contrastive Learning and supervised learning are also involved. These components complement each other, contributing to the superior performance.

#### 3.1 Preliminaries

In our method, we train the model using  $M$  labeled samples and  $N$  unlabeled samples, where  $N \gg M$ . Following (Zhang et al., 2021), we apply the *contextual augmenter* (Shorten et al., 2021) to generate augmented data. Given an unlabeled sample  $\mathbf{x}_i^u$  and a labeled sample  $\mathbf{x}_i^l$ , we denote their initial representations as  $\mathbf{x}_i^{u(0)}$  and  $\mathbf{x}_i^{l(0)}$ , respectively. The corresponding augmented versions are defined as  $\{\mathbf{x}_i^{u(1)}, \mathbf{x}_i^{u(2)}\}$  and  $\{\mathbf{x}_i^{l(1)}, \mathbf{x}_i^{l(2)}\}$ , respectively. During training, mini-batches are constructed from labeled instances  $\mathcal{X} = \{(\mathbf{x}_j^{l(0)}, y_j^l)\}_{j=1}^B$ , and unlabeled instances  $\mathcal{U} = \{(\mathbf{x}_i^{u(0)})_{i=1}^{\mu \cdot B}\}$ . Here,  $B$  is the batch size of labeled data,  $\mu$  is the ratio of unlabeled to labeled examples in each mini-batch, and  $y_j^l$  is the true label corresponding to the cluster  $c \in \{1, \dots, C\}$ . We denote the Encoder as  $f(\cdot)$ , followed by a Classifier network  $g(\cdot)$  and a Projector network  $h(\cdot)$ . For each sample, the probability distributions are defined as

$\mathbf{p}_i \in \mathbb{R}^C = g \circ f(\mathbf{x}_i)$ . The projected representations are defined as  $\mathbf{z}_i \in \mathbb{R}^D = h \circ f(\mathbf{x}_i)$ .

#### 3.2 Interaction-enhanced Optimal Transport

In this section, we will provide a detailed explanation of how sample similarity is introduced to IEOT and leveraged for interaction. By solving this novel OT problem, we can obtain pseudo-labels that effectively integrate both the semantic interactions between individual samples and the global structure from samples to clusters.

Given a batch of original unlabeled texts  $\mathbf{X}_0^u$ , we use the Encoder  $f$  and the Classifier  $g$  to predict probability assignments  $\mathbf{P}_0^u \in \mathbb{R}^{\mu B \times C} = g \circ f(\mathbf{X}_0^u)$ . Then, pseudo-labels can be generated by solving the IEOT problem as follows:

$$\begin{aligned} \min_{\mathbf{Q}, \mathbf{b}} & \langle \mathbf{Q}, \mathbf{M} \rangle + \varepsilon_1 H(\mathbf{Q}) - \varepsilon_2 \Theta(\mathbf{b}) - \varepsilon_3 \langle \mathbf{S}, \mathbf{Q} \mathbf{Q}^T \rangle \\ \text{s.t. } & \mathbf{Q} \mathbf{1} = \mathbf{a}, \mathbf{Q}^T \mathbf{1} = \mathbf{b}, \mathbf{Q} \geq 0, \mathbf{b}^T \mathbf{1} = 1, \end{aligned} \quad (1)$$

where  $\langle \cdot, \cdot \rangle$  is the Frobenius inner product,  $\varepsilon_1, \varepsilon_2$  and  $\varepsilon_3$  are hyperparameters,  $\mathbf{M} = -\log(\mathbf{P}_0^u)$ ,  $\mathbf{S}$  is the cosine similarity matrix of samples. The details of each term in Eq.(1) are as follows:

- $H(\mathbf{Q}) = \langle \mathbf{Q}, \log \mathbf{Q} - \mathbf{1} \rangle$  is the entropy regularization term, which prevents the optimal transport solution from being sparse.
- $\Theta(\mathbf{b}) = \sum_{j=1}^C -b_j \log(b_j)$  is the entropy of the cluster probability  $\mathbf{b}$ , which encourages  $\mathbf{b}$  to approach a uniform distribution. By adjusting the strength of this term, IEOT is suitable for various imbalanced datasets.
- $\langle \mathbf{S}, \mathbf{Q} \mathbf{Q}^T \rangle$  is the semantic regularization, which promotes the transport matrix  $\mathbf{Q}$  to capture semantic similarity interactions between samples. Specifically, this term encourages the transport vector  $\mathbf{Q}_{i:}$  to be similar to  $\mathbf{Q}_{j:}$  when the similarity  $S_{ij}$  is large. In other words, it ensures that semantically similar samples produce similar transport vectors.

IEOT is a novel OT formulation with a complex quadratic semantic regularization term, which traditional OT methods cannot directly solve. Inspired by CSOT, we integrate the Lagrange multiplier algorithm with the generalized conditional gradient (GCG) algorithm to solve the IEOT (Yu et al., 2017). Details of the solution are provided in Appendix A.1. After obtaining transport matrix  $\mathbf{Q}$ ,

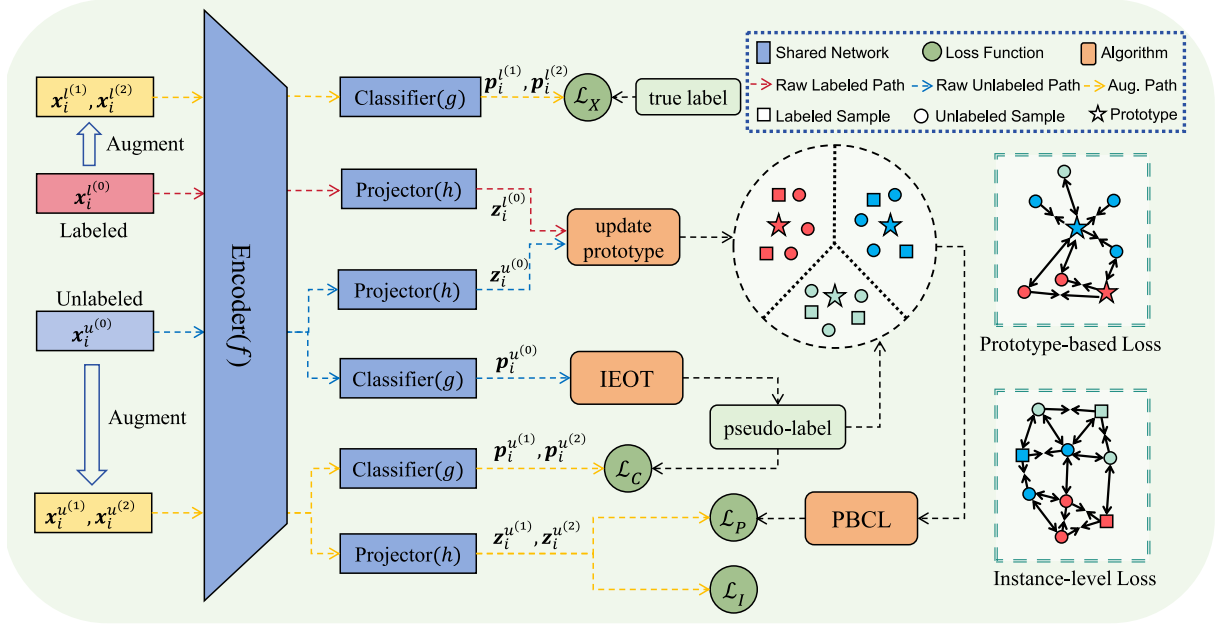


Figure 2: **Method Overview.** IOPC is mainly composed of two core components: Interaction-enhanced Optimal Transport (IEOT) and Prototype-based Contrastive Learning (PBCL).

pseudo-labels  $\hat{Y}^u$  can be generated each element of a batch of pseudo-labels as follows:

$$\hat{Y}_{ij}^u = \begin{cases} 1, & \text{if } j = \arg \max_{j'} Q_{ij'} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

In the IEOT problem, the global minimum-cost matching mechanism introduces global structure information from sample to cluster into the transport matrix  $Q$ , while the term  $\langle S, QQ^T \rangle$  incorporates semantic similarity interaction between samples into the transport matrix  $Q$ , which promotes to generate accurate and reliable pseudo-labels  $\hat{Y}^u$ . Furthermore, we convert the  $i$ -th one-hot pseudo-label in  $\hat{Y}^u$  into the corresponding scalar  $\hat{y}_i^u$ . These reliable pseudo-labels will be further used in the following PBCL, further enhancing the clustering distribution in the feature space.

### 3.3 Prototype-based Contrastive Learning

In this section, we outline the procedure for learning the embedding space defined by the Projector  $h$ . Our goal is to promote well-clustered short text projections by attracting samples to their respective prototypes while distancing them from others. To achieve this goal, we adopt a contrastive objective that utilizes semantic centers as prototypes. Prototypes are computed at the end of each iteration, based on the labeled and confident sam-

ples identified from the previous iteration. Specifically, we utilize a memory bank of size  $\mathcal{O}(2N)$  to record the prediction results  $\hat{y}_i^u$  for each iteration; as well as a reliability indicator for each sample  $\eta_i = \mathbb{1}(\max(\mathbf{p}_i^{u(0)}) \geq \tau)$  denoting if its max prediction exceeds the confidence threshold  $\tau$ .

Subsequently, we update the prototypes  $\mathcal{P} \in \mathbb{R}^{C \times D}$  as the average projections (accumulated over the iteration) of labeled instances and reliable unlabeled instances. Formally, let  $\mathcal{I}_c^l = \{i | \forall \mathbf{x}_i^{l(0)} \in \mathcal{X}, y_i^l = c\}$  be the indices of labeled instances with true cluster  $c$ , and  $\mathcal{I}_c^u = \{i | \forall \mathbf{x}_i^{u(0)} \in \mathcal{U}, \eta_i = 1, \hat{y}_i^u = c\}$  be the indices of the reliable unlabeled samples with hard pseudo-label  $c$ . The normalized prototype for cluster  $c$  can then be obtained as per:

$$\bar{\mathcal{P}}_c = \frac{\sum_{i \in \mathcal{I}_c^u \cup \mathcal{I}_c^l} \mathbf{z}_i}{|\mathcal{I}_c^u| + |\mathcal{I}_c^l|}, \quad \mathcal{P}_c = \frac{\bar{\mathcal{P}}_c}{\|\bar{\mathcal{P}}_c\|_2}. \quad (3)$$

In the following epoch, we minimize the following Prototypical-based Contrastive Learning (PBCL) loss on unlabeled augmented samples:

$$\begin{aligned} \mathcal{L}_P = & -\frac{1}{\mu B} \sum_{i=1}^{\mu B} \log \frac{\exp(\cos(\mathbf{z}_i^{u(1)}, \mathcal{P}_{\hat{y}_i^u})/T_P)}{\sum_{c=1}^C \exp(\cos(\mathbf{z}_i^{u(1)}, \mathcal{P}_c)/T_P)} \\ & -\frac{1}{\mu B} \sum_{i=1}^{\mu B} \log \frac{\exp(\cos(\mathbf{z}_i^{u(2)}, \mathcal{P}_{\hat{y}_i^u})/T_P)}{\sum_{c=1}^C \exp(\cos(\mathbf{z}_i^{u(2)}, \mathcal{P}_c)/T_P)}, \end{aligned} \quad (4)$$



where  $\cos(z_i^{u(1)}, \mathcal{P}_{\hat{y}_i^u})$  means the cosine similarity between  $z_i^{u(1)}$  and the corresponding prototype  $\mathcal{P}_{\hat{y}_i^u}$ , with  $T_P$  meaning the temperature parameter. Consequently, samples from the same category will be more tightly distributed in the feature space.

### 3.4 Instance-level Contrastive Learning

To achieve a more dispersed feature distribution, the instance-level contrastive learning aims to learn representations by bringing positive pairs closer while pushing negative pairs apart. For the  $i$ -th sample, its augmented samples are regarded as a positive pair, while the other  $2N - 2$  pairs are considered negative. The loss function for the  $i$ -th sample is defined as follows:

$$l_i = -\log \frac{\delta(z_i^{u(1)}, z_i^{u(2)})}{\sum_{k=1, k \neq i}^N (\delta(z_i^{u(1)}, z_k^{u(1)}) + \delta(z_i^{u(1)}, z_k^{u(2)}))} - \log \frac{\delta(z_i^{u(2)}, z_i^{u(1)})}{\sum_{k=1, k \neq i}^N (\delta(z_i^{u(2)}, z_k^{u(1)}) + \delta(z_i^{u(2)}, z_k^{u(2)}))}. \quad (5)$$

Here  $\delta(z_i^{u(1)}, z_i^{u(2)}) = \exp(\cos(z_i^{u(1)}, z_i^{u(2)})/T_I)$ ,  $T_I$  is a temperature parameter. The total instance-level contrastive loss is computed as follows:

$$\mathcal{L}_I = \frac{1}{2N} \sum_{i=1}^N l_i. \quad (6)$$

### 3.5 Classification Loss

Using the generated pseudo-labels, we compute the unlabeled loss based on the model's prediction under augmentations, as follows:

$$\mathcal{L}_C = \frac{1}{\mu B} \sum_{i=1}^{\mu B} (\text{CE}(\hat{y}_i^u, \mathbf{p}_i^{u(1)}) + \text{CE}(\hat{y}_i^u, \mathbf{p}_i^{u(2)})), \quad (7)$$

where CE denotes cross-entropy. Also, we apply a supervised classification loss over the labeled data:

$$\mathcal{L}_X = \frac{1}{B} \sum_{i=1}^B (\text{CE}(y_i^l, \mathbf{p}_i^{l(1)}) + \text{CE}(y_i^l, \mathbf{p}_i^{l(2)})), \quad (8)$$

Notably, Eq.(7) and Eq.(8) are acted on both two augmented versions.

### 3.6 Final Objective

We design a two-stage training procedure for IOPC. The first stage primarily aims to obtain a good initial feature space, while the second stage focuses

on optimizing the distribution using all the algorithms mentioned above. The overall loss function of the model is as follows:

$$\mathcal{L} = \begin{cases} \mathcal{L}_{S_1} = \mathcal{L}_X + \mathcal{L}_C + \mathcal{L}_I & \text{if } iter < 1000 \\ \mathcal{L}_{S_2} = \lambda(\mathcal{L}_{S_1}) + \mathcal{L}_P & \text{if } iter \geq 1000, \end{cases} \quad (9)$$

where  $iter$  is the number of training iterations,  $\lambda$  is a balancing hyperparameter. By integrating the above components, the model learns a high-quality feature space distribution, leading to more accurate and stable clustering results. Algorithm 3 in Appendix B.1 describes the training process of IOPC.

## 4 Experiments

### 4.1 Datasets

We conducted experiments using eight benchmark datasets: **AgNews**, **StackOverflow**, **Biomedical**, **SearchSnippets**, **GoogleNews-TS**, **GoogleNews-T**, **GoogleNews-S**, and **Tweet**. A summary of the key characteristics and detailed information of these datasets are provided in Table 1 and Appendix B.2, respectively.

Datasets	S	N	L	R
AgNews	8000	4	23	1
SearchSnippets	12340	8	18	7
StackOverflow	20000	20	8	1
Biomedical	20000	20	13	1
GoogleNews-TS	11109	152	8	143
GoogleNews-T	11109	152	6	143
GoogleNews-S	11109	152	22	143
Tweet	2472	89	22	249

Table 1: **Key Information of Datasets.** "S" represent the dataset size; "N" is the number of categories; "L" is the average sentence length; "R" is the size ratio of the largest to the smallest category.

### 4.2 Experiment Settings

We implement our model using PyTorch (Paszke et al., 2019) and employ *bge-base-en-v1.5* in the Sentence Transformers library as the Encoder (Chen et al., 2024). Under our few-shot definition, we use 1% of the samples as labeled samples if  $S/N > 1\%$  according to Table 1, otherwise we use only 1 sample per dataset as labeled samples. All parameters of our model are optimized using the Adam optimizer (Kingma, 2014). The learning rate of the Encoder is  $5 \times 10^{-6}$ , while the other networks is  $5 \times 10^{-4}$ . We use Accuracy (ACC) and Normalized Mutual Information (NMI) to evaluate

	AgNews		SearchSnippets		StackOverflow		Biomedical	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
SCCL	83.10	61.96	79.90	63.78	70.83	69.21	42.49	39.16
RSTC	84.24	62.45	80.10	69.74	83.30	74.11	48.40	40.12
MIST	89.47	70.25	76.72	67.69	79.65	78.59	39.15	34.66
STSPL-SSC	<u>89.92</u>	<u>71.66</u>	81.04	65.46	86.74	<u>82.54</u>	47.43	42.49
COTC	87.56	67.09	<u>90.32</u>	<u>77.09</u>	<u>87.78</u>	79.19	<u>53.20</u>	<u>46.09</u>
BGE-M3	87.89	66.67	75.59	60.7	84.66	82.21	51.25	46.05
<b>IOPC</b>	<b>90.28</b>	<b>72.22</b>	<b>90.44</b>	<b>77.15</b>	<b>90.38</b>	<b>82.74</b>	<b>60.54</b>	<b>48.81</b>
<b>Improvement</b>	<b>+0.36</b>	<b>+0.56</b>	<b>+0.12</b>	<b>+0.06</b>	<b>+2.6</b>	<b>+0.20</b>	<b>+7.34</b>	<b>+2.72</b>
	GoogleNews-TS		GoogleNews-T		GoogleNews-S		Tweet	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
SCCL	82.51	93.01	69.01	85.10	73.44	87.98	73.10	86.66
RSTC	83.27	93.15	72.27	87.39	79.32	89.40	75.20	87.35
MIST	<u>90.63</u>	<b>96.42</b>	78.80	89.31	82.14	90.86	<u>91.75</u>	<b>95.12</b>
STSPL-SSC	84.41	94.32	81.01	91.11	82.30	91.18	79.59	88.02
COTC	90.50	<u>96.33</u>	<u>83.53</u>	<u>92.07</u>	<u>86.10</u>	<b>93.49</b>	91.33	<u>95.09</u>
BGE-M3	72.97	91.81	68.28	87.52	69.89	89.01	64.64	87.42
<b>IOPC</b>	<b>92.92</b>	95.90	<b>87.71</b>	<b>92.39</b>	<b>87.64</b>	<u>92.79</u>	<b>92.11</b>	94.63
<b>Improvement</b>	<b>+2.29</b>	<b>-0.52</b>	<b>+4.18</b>	<b>+0.32</b>	<b>+1.54</b>	<b>-0.7</b>	<b>+0.36</b>	<b>-0.49</b>

Table 2: **Experimental Results.** Here we present the clustering performance of our method IOPC and the baselines on eight benchmarks. The results for the baselines are quoted from (Zheng et al., 2023; Li et al., 2024; Kamthawee et al., 2024; Nie et al., 2024). We bold the **best result**, underline the runner-up.

the model. Definitions of the metrics and detailed settings are in Appendix B.3 and Appendix B.4.

### 4.3 Baselines

We compare **IOPC** with several latest short text clustering approaches. **SCCL** (Zhang et al., 2021) employs contrastive learning to refine representations and obtains the clustering results using the DEC algorithm (Xie et al., 2016). **RSTC** (Zheng et al., 2023) constructs pseudo-labels using adaptive optimal transport to assist the model in training neural networks for clustering. **MIST** (Kamthawee et al., 2024) enhances clustering by maximizing the mutual information between representations at both the sequence and token levels. **STSPL-SSC** (Nie et al., 2024) extends RSTC by incorporating additional labeled data and leveraging the information from these labels to guide the effectiveness of pseudo-labels. **COTC** (Li et al., 2024) introduces a Co-Training Clustering framework that effectively combines BERT and TFIDF features to generate a high-quality feature space for clustering.

Additionally, to measure the efficiency of the Encoder, we include the **BGE-M3** experiment, which applies k-means directly to the output of the BGE-

M3 model. Further analysis of this Encoder on other baselines are illustrated in the Appendix C.3.

### 4.4 Main Results

The clustering results for both baseline models and IOPC are summarized in Table 2. From the results, we can find that: (1) The traditional contrastive learning method **SCCL** and the **RSTC** method with the introduction of OT, due to the complexity of the datasets, did not yield good results. (2) Directly incorporating k-means in **BGE-M3** cannot achieve good clustering results. (3) **MIST** and **COTC** allow the model to learn more features, and thus performed second only to IOPC on some datasets. However, they still struggled to address the challenges posed by complex datasets. (4) **STSPL-SSC**, by introducing semi-supervised learning, demonstrated good performance; nevertheless, the information it could learn still fell short of IOPC, so did its performance. (5) Obviously, **IOPC** consistently outperforms previous methods across all datasets. Notably, IOPC achieves superior clustering accuracy, particularly on more challenging datasets such as Biomedical, GoogleNews-T, and StackOverflow. The two components in

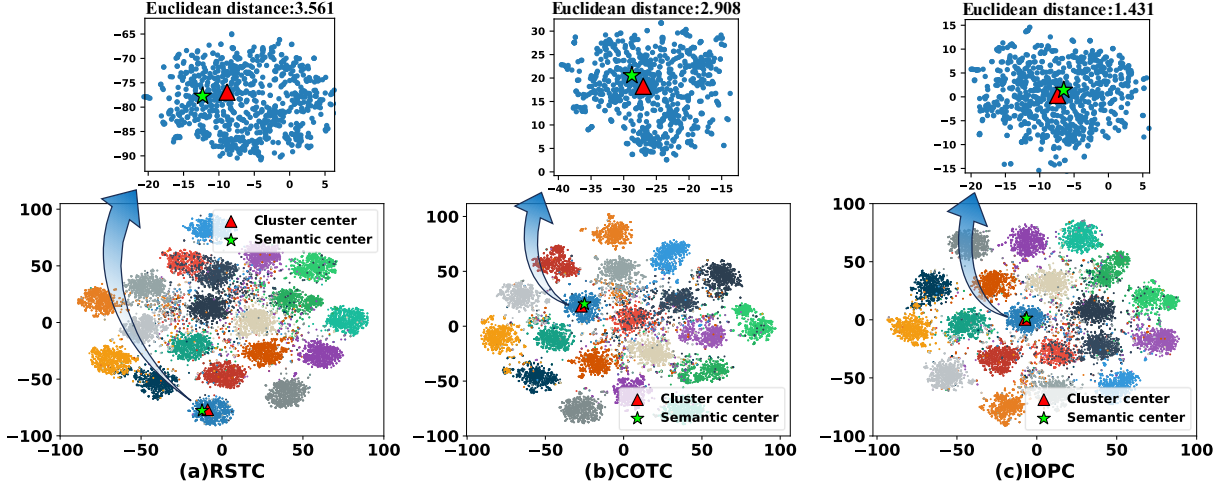


Figure 3: **Comparison of the Alignment Between Semantic Centers and Cluster Centers.** The euclidean distance is computed based on the T-SNE mapping. Each color indicates a truth category.

IOPC cooperate with each other to extract such scarce information, achieving a more regular and stable distribution in the feature space, which is essential for achieving such outstanding results. In the following sections, numerous experiments such as visualization will be presented to further validate the accuracy and stability of our model.

#### 4.5 Semantic Alignment Visualization

We use T-SNE visualization and Euclidean distances to verify whether IOPC achieves semantic alignment. Specifically, we used the "Matlab" category from the StackOverflow dataset as an example, with the semantic center generated from a sentence most relevant to the keywords from this category (i.e. "Matlab functions matrices visualization programming scripts optimization"). The results are shown in Figure 3, compared to other models, the cluster center of our model are closest to the semantic center. It reveals that our method achieves the best alignment of cluster centers with the semantic centers.

Furthermore, we can observe that the feature space distribution obtained by IOPC is more consistent and compact. A more detailed comparison of the representation visualizations is provided in Appendix C.2.

#### 4.6 The Comparison of Model Stability

To validate the stability of our model, we used multiple different random seeds to observe variations in model performance. Specifically, we conducted experiments on the AgNews and SearchSnippets datasets, with random seeds ranging from

0 to 10. To ensure a fair comparison, all experiments uniformly use BGE-M3 as Encoder. The results are shown in Figure 4. From it, we can find that: (1) RSTC demonstrates high stability but performs poorly on the imbalanced SearchSnippets dataset. (2) COTC exhibits lower stability. (3) IOPC achieves the highest performance while maintaining strong stability, demonstrating the robustness and generalizability of our model.

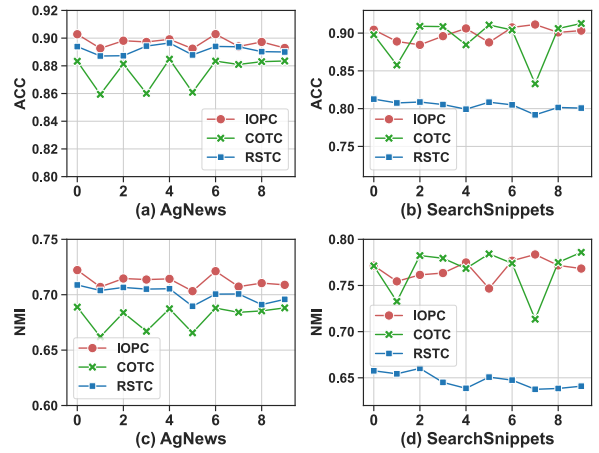


Figure 4: **Comparison of Stability.** The x-axis representing the random seeds we used.

#### 4.7 Ablation Study

To demonstrate that each proposed modification in IOPC contributes to more accurate clustering, we conducted ablation experiments on eight datasets, and the experimental results are presented in Table 3. The experimental results demonstrate that the model performance significantly decreases regardless of which module we remove from our proposed approach. When PBCL is removed, relying solely

Losses	Agn	Sea	Sta	Bio	GN-TS	GN-T	GN-S	Twe	$\delta$
$-(\mathcal{L}_C, \mathcal{L}_P)$	86.50	81.70	86.74	49.40	81.13	64.79	73.04	73.21	-11.94
$-\mathcal{L}_C$	87.41	84.24	88.10	53.51	82.77	67.23	74.86	75.32	-9.82
$-\mathcal{L}_P$	88.79	87.51	89.33	58.17	91.47	86.42	86.48	90.41	-1.68
<b>IOPC</b>	90.28	90.44	90.38	60.54	92.92	87.71	87.64	92.11	0

Table 3: **Ablation Results.**  $-\mathcal{L}_*$  denotes the respective loss is not applied.  $\delta$  is the average improvement over IOPC.

Labeled count	Agn	Sea	Sta	Bio	GN-TS	GN-T	GN-S	Twe
1 or 1%	90.28	90.44	90.38	60.54	92.92	87.71	87.64	92.11
2 or 2%	90.41	91.13	90.83	63.51	94.21	89.1	90.86	94.7
5 or 5%	91.13	92.35	91.22	69.43	95.02	90.35	91.17	95.23
10 or 10%	91.65	93.25	91.96	73.41	96.25	93.09	92.84	98.46

Table 4: **The Impact of Varying the Number of Labeled Samples.** Note that, when  $(S/N) \leq 1\%$ , if the required labeled samples for a class exceed its available samples, the available number of samples in that class is used instead.

on IEOT to generate pseudo-labels fails to optimize the distribution in the feature space. On the other hand, when IEOT is removed, PBCL cannot utilize reliable pseudo-labels, causing the failure in learning the correct information. Only when each part of the model collaborates with the others can the best performance be achieved.

#### 4.8 The Impact of Labeled Data Quantity

We conduct experiments by varying the number of labeled samples to 1 or 1%, 2 or 2%, 5 or 5%, 10 or 10%, where "1 or 1%" means that: we use 1% of the samples as labeled if  $(S/N) > 1\%$  according to Table 1, and we use one sample per category as labeled if  $(S/N) \leq 1\%$ . The results are presented in Table 4. We can observe that the performance increases with the number of labeled samples. In few-shot settings, our model already achieves state-of-the-art results, and as more labeled data is collected, the model's performance continues to improve. This demonstrates that our model can effectively be applied in real-world scenarios. Finally, we construct the labeled data using the "1 or 1%" setting, which offers the highest cost-effectiveness.

#### 4.9 In-depth Analysis

In addition to the experiments mentioned above, we conducted several supplementary experiments to further verify the capabilities of IOPC.

(1) We recorded how the number of predicted clusters are changing over iterations in Appendix C.1, showing that our model can effectively combat clustering degeneracy. (2) Since each baseline model uses a different Encoder, we converted baseline models to the same Encoder (BGE-M3) for comparison. The results provided in the Appendix

C.3, it can be observed that, regardless of whether the Encoder is the same or not, our model outperforms all other models. (3) Due to the current scarcity of semi-supervised methods in the field of short text clustering, we incorporated labeled data into recent high-performance models in the training process. As can be seen from the Appendix C.4, few-shot scenario will not directly enhance the performance of the previous baselines, and IOPC still outperforms these models comprehensively. (4) We conducted hyperparameter analysis experiments including  $\varepsilon_1, \varepsilon_2, \varepsilon_3$  and  $\lambda$ , and analyzed the impact of these hyperparameters on IOPC which can be found in the Appendix E. (5) We recorded the computation budget with previous models, as shown in Appendix D. Our model strikes a balance between performance and efficiency, making it the most cost-effective solution.

## 5 Conclusion

This paper presents a novel approach, **IOPC**, for few-shot short text clustering, which combines Interaction-enhanced Optimal Transport (**IEOT**) and Prototype-based Contrastive Learning (**PBCL**). The former significantly improved the accuracy of pseudo-labels by exploiting the interaction between samples, while the latter using semantic centers to align the cluster centers by leveraging labeled and confident pseudo-labeled data to construct prototypes. Thereby, achieving high-quality feature space distribution. Extensive experiments demonstrate that our model consistently outperforms existing state-of-the-art techniques, showing significant improvements in clustering accuracy and stability.



## 6 Limitations

Despite the promising results, there are some limitations to our method. (1) The performance slightly depends on the quality and representativeness of the labeled data. So the future work will focus on how to derive labeled data in a cost-effective way like using LLMs. (2) The pseudo-labeling process, while effective, can still introduce errors, particularly in noisy or ambiguous data. Therefore, exploring a method for generating more accurate pseudo-labels is also a key focus in the future.

## References

Sepideh Bazzaz Abkenar, Mostafa Haghi Kashani, Mohammad Akbari, and Ebrahim Mahdipour. 2023. Learning textual features for twitter spam detection: A systematic literature review. volume 228, page 120366. Elsevier.

Eric Arazo, Diego Ortego, Paul Albert, Noel E O'Connor, and Kevin McGuinness. 2020. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International joint conference on neural networks (IJCNN)*, pages 1–8. IEEE.

Liang Bai, Jiye Liang, Chuangyin Dang, and Fuyuan Cao. 2012. A cluster centers initialization method for clustering categorical data. *Expert Systems with Applications*, 39(9):8022–8029.

Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. [Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation](#). In *Annual Meeting of the Association for Computational Linguistics*.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR.

W Dong-DongChen and ZH WeiGao. 2018. Tri-net for semi-supervised deep learning. In *Proceedings of twenty-seventh international joint conference on artificial intelligence*, pages 2014–2020.

Enrico Fini, Pietro Astolfi, Karteek Alahari, Xavier Alameda-Pineda, Julien Mairal, Moin Nabi, and Elisa Ricci. 2023. Semi-supervised learning made simple with self-supervised clustering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3187–3197.

Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al.

2020. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284.

Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE.

Krissanee Kamthawee, Can Udomcharoenchaikit, and Sarana Nutanong. 2024. Mist: mutual information maximization for short text clustering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11309–11324.

Diederik P Kingma. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Mohammad Amin Kuhail, Nazik Alturki, Salwa Alramlawi, and Kholood Alhejori. 2023. Interacting with educational chatbots: A systematic review. *Education and Information Technologies*, 28(1):973–1018.

Dong-Hyun Lee et al. 2013. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, page 896. Atlanta.

Zetong Li, Qinliang Su, Shijing Si, and Jianxing Yu. 2024. Leveraging BERT and TFIDF Features for Short Text Clustering via Alignment-Promoting Co-Training. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 14897–14913.

Belal Abdullah Hezam Murshed, Suresha Mallappa, Jemal Abawajy, Mufeed Ahmed Naji Saif, Hasib Daowd Esmail Al-Ariki, and Hudhaifa Mohammed Abdulwahab. 2023. Short text topic modelling approaches in the context of big data: taxonomy, survey, and analysis. volume 56, pages 5133–5260. Springer.

Islam Nassar, Samitha Herath, Ehsan Abbasnejad, Wray Buntine, and Gholamreza Haffari. 2021. All labels are not created equal: Enhancing semi-supervision via label grouping and co-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7241–7250.

Wenhua Nie, Lin Deng, Chang-Bo Liu, JialingWei JialingWei, Ruitong Han, and Haoran Zheng. 2024. STSPL-SSC: Semi-Supervised Few-Shot Short Text Clustering with Semantic text similarity Optimized Pseudo-Labels. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 12174–12185, Bangkok, Thailand. Association for Computational Linguistics.

Christos H Papadimitriou and Kenneth Steiglitz. 1998. *Combinatorial optimization: algorithms and complexity*. Courier Corporation.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.

Xuan-Hieu Phan, Le-Minh Nguyen, and Susumu Horiguchi. 2008. Learning to classify short and sparse text & web with hidden topics from large-scale data collections. In *Proceedings of the 17th international conference on World Wide Web*, pages 91–100.

Md Rashadul Hasan Rakib, Norbert Zeh, Magdalena Jankowska, and Evangelos Milios. 2020. Enhancement of short text clustering by iterative classification. In *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 25*, pages 105–117. Springer.

Connor Shorten, Taghi M Khoshgoftaar, and Borko Furht. 2021. Text data augmentation for deep learning. *Journal of big Data*, 8(1):101.

Phani Teja Nallamothu and Mohd Shais Khan. 2023. Machine learning for spam detection. *Asian Journal of Advances in Research*, 6(1):167–179.

Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International conference on machine learning*, pages 9929–9939. PMLR.

Junyuan Xie, Ross Girshick, and Ali Farhadi. 2016. Unsupervised deep embedding for clustering analysis. In *International conference on machine learning*, pages 478–487. PMLR.

Jiaming Xu, Bo Xu, Peng Wang, Suncong Zheng, Guanhua Tian, and Jun Zhao. 2017. Self-taught convolutional neural networks for short text clustering. *Neural Networks*, 88:22–31.

Jianhua Yin and Jianyong Wang. 2016. A model-based approach for text clustering with outlier detection. In *2016 IEEE 32nd International Conference on Data Engineering (ICDE)*, pages 625–636. IEEE.

Yaoliang Yu, Xinhua Zhang, and Dale Schuurmans. 2017. Generalized conditional gradient for sparse estimation. *Journal of Machine Learning Research*, 18(144):1–46.

Dejiao Zhang, Feng Nan, Xiaokai Wei, Shangwen Li, Henghui Zhu, Kathleen McKeown, Ramesh Nallapati, Andrew Arnold, and Bing Xiang. 2021. Supporting clustering with contrastive learning. *arXiv preprint arXiv:2103.12953*.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.

Xiaolin Zheng, Mengling Hu, Weiming Liu, Chaochao Chen, and Xinting Liao. 2023. Robust Representation Learning with Reliable Pseudo-labels Generation via Self-Adaptive Optimal Transport for Short Text Clustering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10493–10507.

## A Hyper-efficient Solution for IEOT

### A.1 Formulation of the Solution

As mentioned in Section 3.2, the IEOT problem is formulated as:

$$\begin{aligned} \min_{\mathbf{Q}, \mathbf{b}} & \langle \mathbf{Q}, \mathbf{M} \rangle + \varepsilon_1 H(\mathbf{Q}) - \varepsilon_2 \Theta(\mathbf{b}) - \varepsilon_3 \langle \mathbf{S}, \mathbf{Q} \mathbf{Q}^T \rangle \\ \text{s.t.} & \quad \mathbf{Q} \mathbf{1} = \mathbf{a}, \quad \mathbf{Q}^T \mathbf{1} = \mathbf{b}, \quad \mathbf{Q} \geq 0, \quad \mathbf{b}^T \mathbf{1} = 1, \end{aligned} \quad (10)$$

where  $\langle \cdot, \cdot \rangle$  represents the Frobenius inner product,  $\varepsilon_1$  and  $\varepsilon_2$  are balancing hyperparameters,  $H(\mathbf{Q}) = \langle \mathbf{Q}, \log(\mathbf{Q}) - \mathbf{1} \rangle$ , and  $\Theta(\mathbf{b}) = \sum_{j=1}^C -b_j \log(b_j)$  is the entropy of the cluster probability assignments  $\mathbf{b}$ .

The IEOT incorporates a complex quadratic semantic regularization term, which cannot be solved directly using traditional OT methods. To address this OT problem, inspired by CSOT, we propose integrating the Lagrange multiplier algorithm (Zheng et al., 2023) into the generalized conditional gradient (GCG) algorithm (Yu et al., 2017) to solve the IEOT problem. Specifically, we first utilize the GCG algorithm linearize the complex quadratic term  $\langle \mathbf{S}, \mathbf{Q} \mathbf{Q}^T \rangle$ , then we employ the Lagrange multiplier algorithm to solve it.

To better explain the linearization of the semantic constraint term  $\langle \mathbf{S}, \mathbf{Q} \mathbf{Q}^T \rangle$ , we first define:

$$f(\mathbf{Q}) = \langle \mathbf{Q}, \mathbf{M} \rangle - \varepsilon_3 \langle \mathbf{S}, \mathbf{Q} \mathbf{Q}^T \rangle, \quad (11)$$

the algorithm GCG is an iterative optimization algorithm. For the  $i$ -th iteration, given the previously computed transport matrix  $\mathbf{Q}_{i-1}$ , the objective function  $f(\mathbf{Q})$  is expanded using a Taylor series and only retained the linear term:

$$\begin{aligned} f_{\text{lin}}(\mathbf{Q}) &= f(\mathbf{Q}_{i-1}) + \langle f'(\mathbf{Q}_{i-1}), \mathbf{Q} - \mathbf{Q}_{i-1} \rangle \\ &= \langle \mathbf{Q}, f'(\mathbf{Q}_{i-1}) \rangle + f(\mathbf{Q}_{i-1}) - \langle \mathbf{Q}_{i-1}, f'(\mathbf{Q}_{i-1}) \rangle, \end{aligned} \quad (12)$$

since  $\mathbf{Q}_{i-1}$  has already been computed, the term  $f(\mathbf{Q}_{i-1}) - \langle \mathbf{Q}_{i-1}, f'(\mathbf{Q}_{i-1}) \rangle$  is a constant, we use the letter  $C$  to represent it. Therefore, the Eq.(11)

can be approximate to:

$$\begin{aligned} f(\mathbf{Q}) &\approx f_{\text{in}}(\mathbf{Q}) \\ &= \langle \mathbf{Q}, f'(\mathbf{Q}_{i-1}) \rangle + f(\mathbf{Q}_{i-1}) - \langle \mathbf{Q}_{i-1}, f'(\mathbf{Q}_{i-1}) \rangle \\ &= \langle \mathbf{Q}, f'(\mathbf{Q}_{i-1}) \rangle + C. \end{aligned} \quad (13)$$

Substituting Eq.(13) into Eq.(10), the IEOT problem can be approximately simplified as:

$$\begin{aligned} \min_{\mathbf{Q}, \mathbf{b}} & \langle \mathbf{Q}, f'(\mathbf{Q}_{i-1}) \rangle + \varepsilon_1 H(\mathbf{Q}) - \varepsilon_2 \Theta(\mathbf{b}) \\ \text{s.t. } & \mathbf{Q}\mathbf{1} = \mathbf{a}, \mathbf{Q}^T \mathbf{1} = \mathbf{b}, \mathbf{Q} \geq 0, \mathbf{b}^T \mathbf{1} = 1, \end{aligned} \quad (14)$$

where the formulation of  $f'(\mathbf{Q}_{i-1})$  is provided in Appendix A.2. The pseudo-code of the GCG algorithm is presented in Algorithm 2.

**Algorithm 1** Generalized Conditional Gradient Algorithm for IEOT with Quadratic Constraints

**Input:** Probability matrix  $\mathbf{P}^{(0)}$ ; marginal constraints  $\mathbf{a}$ ; semantic similarity matrix  $\mathbf{S}$ ; constraints weights  $\varepsilon_1, \varepsilon_2$  and  $\varepsilon_3$ .

**Output:** Transport matrix  $\mathbf{Q}$ .

Initialize  $\mathbf{b}_0$  randomly and perform normalization so that  $\mathbf{b}_0^T \mathbf{1} = 1$

Initialize  $\mathbf{Q}_0 = \mathbf{a}\mathbf{b}_0^T$ .

**for**  $i = 1$  to  $T_1$  **do**

$$f'(\mathbf{Q}_{i-1}) = \mathbf{M} - \varepsilon_3(\mathbf{S} + \mathbf{S}^T)\mathbf{Q}_{i-1}.$$

$$\tilde{\mathbf{Q}}_i, \mathbf{b}_i = \arg \min_{\mathbf{Q}, \mathbf{b}} \langle \mathbf{Q}, f'(\mathbf{Q}_{i-1}) \rangle + \varepsilon_1 H(\mathbf{Q}) + \varepsilon_2 \Theta(\mathbf{b}).$$

Choose  $\alpha_i \in [0, 1]$  so that it satisfies the Armijo rule.

$$\mathbf{Q}_i = (1 - \alpha_i)\mathbf{Q}_{i-1} + \alpha_i \tilde{\mathbf{Q}}_i.$$

**end for**

Then, we adopt the Lagrangian multiplier algorithm to solve Eq.(14):

$$\begin{aligned} \min_{\mathbf{Q}, \mathbf{b}} & \langle \mathbf{Q}, \mathbf{M}' \rangle + \varepsilon_1 H(\mathbf{Q}) - \varepsilon_2 \Theta(\mathbf{b}) - \mathbf{f}^T(\mathbf{Q}\mathbf{1} - \mathbf{a}) \\ & - \mathbf{g}^T(\mathbf{Q}^T \mathbf{1} - \mathbf{b}) - h(\mathbf{b}^T \mathbf{1} - 1), \end{aligned} \quad (15)$$

where  $\mathbf{f}, \mathbf{g}$  and  $h$  are all Lagrangian multipliers,  $\mathbf{M}' = f'(\mathbf{Q}_{i-1})$ . Taking the partial derivative of Eq.(15) with respect to  $\mathbf{Q}$ , we can obtain:

$$Q_{ij} = \exp\left(\frac{f_i + g_j - M'_{ij}}{\varepsilon_1}\right) > 0. \quad (16)$$

Eq.(16) is a function of each element in  $\mathbf{f}$  and  $\mathbf{g}$ . Next, we first fix  $\mathbf{b}$ , and update  $f_i$  and  $g_j$ . Due to

the fact that  $\mathbf{Q}\mathbf{1} = \mathbf{a}$ , we can get:

$$\begin{aligned} \sum_{j=1}^C Q_{ij} &= \sum_{j=1}^C \exp\left(\frac{f_i + g_j - M'_{ij}}{\varepsilon_1}\right) \\ &= \exp\left(\frac{f_i}{\varepsilon_1}\right) \sum_{j=1}^C \exp\left(\frac{g_j - M'_{ij}}{\varepsilon_1}\right) \\ &= a_i, \end{aligned} \quad (17)$$

where  $C$  represents the number of clusters in the dataset. Further, we can obtain:

$$\exp\left(\frac{f_i}{\varepsilon_1}\right) = \frac{a_i}{\sum_{j=1}^C \exp\left(\frac{g_j - M'_{ij}}{\varepsilon_1}\right)}. \quad (18)$$

taking the logarithm of both sides and multiplying by  $\varepsilon_1$ , we can obtain:

$$f_i = \varepsilon_1 \ln a_i - \varepsilon_1 \ln \sum_{j=1}^C \exp\left(\frac{g_j - M'_{ij}}{\varepsilon_1}\right). \quad (19)$$

Similar to the above derivation, from  $\mathbf{Q}^T \mathbf{1} = \mathbf{b}$ , we can obtain:

$$g_j = \varepsilon_1 \ln b_j - \varepsilon_1 \ln \sum_{i=1}^{\mu B} \exp\left(\frac{f_i - M'_{ij}}{\varepsilon_1}\right). \quad (20)$$

We can observe that  $g_j$  is an unknown variable in Eq.(19), while  $f_i$  is an unknown variable in Eq.(20). Since  $f_i$  and  $g_j$  are functions of each other, making it infeasible to directly solve for their exact values. Thus, we employ an iterative approach to update and work out it.

Then, we fix  $\mathbf{f}$  and  $\mathbf{g}$ , and update  $\mathbf{b}$ . Specifically, take the partial derivative of the optimization problem Eq.(15) on the variable  $\mathbf{b}$ , we can obtain:

$$\varepsilon_2(\log(b_j) + 1) + g_j - h = 0, \quad (21)$$

by solving formula Eq.(21), we can get:

$$b_j(h) = \exp\left(\frac{h - g_j - \varepsilon_2}{\varepsilon_2}\right). \quad (22)$$

Taking Eq.(22) back to the original constraint  $\mathbf{b}^T \mathbf{1} = 1$ , the formula is defined as below:

$$(\mathbf{b}(h))^T \mathbf{1} = \sum_{j=1}^C \exp\left(\frac{h - g_j - \varepsilon_2}{\varepsilon_2}\right) = 1, \quad (23)$$

by extracting the scalar part, we obtain:

$$\exp\left(\frac{h}{\varepsilon_2}\right) \sum_{j=1}^C \exp\left(\frac{-g_j - \varepsilon_2}{\varepsilon_2}\right) = 1, \quad (24)$$

by solving Eq.(24), we can get:

$$h = -\varepsilon_2 \log \left( \sum_{j=1}^C \exp\left(\frac{-g_j - \varepsilon_2}{\varepsilon_2}\right) \right), \quad (25)$$

where  $h$  is the root of Eq.(23), Then, we can obtain  $\mathbf{b}$  by Eq.(22).

Overall, through iteratively updating the Eq.(19), (20) and (22), we can get the transport matrix  $\mathbf{Q}$  on Eq.(16). We show the iterative optimization process for solving Eq.(14) using the Lagrange multiplier algorithm in Algorithm 2.

---

**Algorithm 2** The optimization scheme of IEOT

---

**Input:** The cost distance matrix  $\mathbf{M}'$ ; cluster probability  $b$ ;

**Output:** The transport matrix  $\mathbf{Q}$

**Procedure:**

Initialize  $\mathbf{f}$  and  $\mathbf{g}$  randomly.

Initialize  $h = 1$ .

**for**  $i=1$  to  $T_2$  **do**

Fix  $\mathbf{b}$ , update  $\mathbf{f}$  and  $\mathbf{g}$  by Eq.(19) and (20), respectively.

Fix  $\mathbf{f}$  and  $\mathbf{g}$ , update  $\mathbf{b}$  by Eq.(22) and (25).

**end for**

Calculate  $\mathbf{Q}$  in Eq.(16).

---

## A.2 Derivative Complement

The calculation of  $f'(Q_{i-1})$  in Eq.(14) is as follows:

$$\begin{aligned} \frac{\partial f(Q_{i-1})}{\partial Q} &= \frac{\partial \langle Q, M \rangle - \varepsilon_3 \langle S, QQ^T \rangle}{\partial Q} \\ &= M - \varepsilon_3 \frac{\partial \langle S, QQ^T \rangle}{\partial Q}, \end{aligned} \quad (26)$$

where:

$$\begin{aligned} \frac{\partial \langle S, QQ^T \rangle}{\partial Q} &= \frac{\partial \text{tr}(S^T QQ^T)}{\partial Q} = \frac{\partial \text{tr}((S^T QQ^T)^T)}{\partial Q} \\ &= \frac{\partial \text{tr}(QQ^T S)}{\partial Q} = \frac{\partial \text{tr}(Q^T S Q)}{\partial Q} \\ &= (S + S^T)Q, \end{aligned} \quad (27)$$

therefore, the final computation result is as follows:

$$\frac{\partial f(Q_{i-1})}{\partial Q} = M - \varepsilon_3 (S + S^T)Q. \quad (28)$$

## B Supplementary Details

### B.1 Pseudocode of IOPC

We present the pseudocode of IOPC's training process for an iteration, as shown in Algorithm 3.

---

**Algorithm 3** Pseudocode for an iteration of IOPC

---

**Input:** Encoder  $f$ ; Classifier  $g$ ; Projector  $h$ ; Mini-batch labeled data  $\{X^{l(0)}, Y^l\}$ ; Mini-batch unlabeled data  $X^{u(0)}$ ; current iteration  $iter$ .

**Output:** Updated parameters

# generate augmented samples

$X^{l(1)}, X^{l(2)} \leftarrow \text{textual augmenter}(X^{l(0)})$

$X^{u(1)}, X^{u(2)} \leftarrow \text{textual augmenter}(X^{u(0)})$

# forward texts and obtain  $\mathbf{P}$  and  $\mathbf{Z}$

$\mathbf{P}^{l(1)}, \mathbf{P}^{l(2)}, \mathbf{P}^{u(0)}, \mathbf{P}^{u(1)}, \mathbf{P}^{u(2)} \leftarrow f(g(\sim))$

$\mathbf{Z}^{l(0)}, \mathbf{Z}^{u(0)}, \mathbf{Z}^{u(1)}, \mathbf{Z}^{u(2)} \leftarrow f(h(\sim))$

# produce pseudo-label via IEOT

$\hat{Y}^u \leftarrow \text{IEOT}(\mathbf{P}^{u(0)})$  # Eq.(2)

# accumulate prototypes

$\eta \leftarrow \mathbb{1}(\max(\mathbf{P}^{u(0)}) \geq \tau)$

$\bar{\mathcal{P}} \leftarrow \text{accum\_prototypes}(\mathbf{Z}^{l(0)}, \mathbf{Z}^{u(0)}, \eta)$  # Eq.(3)

$\mathcal{P} \leftarrow \text{update\_prototypes}(\bar{\mathcal{P}})$  # Eq.(3)

# calculate the loss function

$\mathcal{L}_I \leftarrow \text{calculate loss}(\mathbf{Z}^{u(1)}, \mathbf{Z}^{u(2)})$  # Eq.(6)

$\mathcal{L}_X \leftarrow \text{calculate loss}(\mathbf{P}^{l(1)}, \mathbf{P}^{l(2)}, \mathbf{Y}^l)$  # Eq.(8)

$\mathcal{L}_C \leftarrow \text{calculate loss}(\mathbf{P}^{u(1)}, \mathbf{P}^{u(2)}, \hat{\mathbf{Y}}^u)$  # Eq.(7)

$\mathcal{L} \leftarrow \mathcal{L}_I + \mathcal{L}_X + \mathcal{L}_C$  # Eq.(9)

**if**  $iter < 1000$  **then**

$\mathcal{L}_P \leftarrow \text{calculate loss}(\mathbf{Z}^{u(1)}, \mathbf{Z}^{u(2)}, \hat{\mathbf{Y}}^u, \mathcal{P})$  # Eq.(4)

$\mathcal{L} \leftarrow \mathcal{L} + \mathcal{L}_P$  # Eq.(9)

**end if**

# update parameters

$\text{back\_propagation}(\mathcal{L})$

---

### B.2 Datasets

We evaluated the performance of IOPC on eight benchmark datasets, which cover a wide range of text sources, including news headlines and social media content. These diverse sets enable a thorough evaluation of the model across various domains. Based on the degree of imbalance, **AgNews**, **StackOverflow**, and **Biomedical** are classified as balanced datasets, while **SearchSnippets** is categorized as a slightly imbalanced dataset. In contrast, **GoogleNews-TS**, **GoogleNews-T**, **GoogleNews-S**, and **Tweet** are considered as severely imbalanced datasets. The brief descriptions are provided below:

- **AgNews:** Sourced from AG's news corpus (Zhang et al., 2015), this dataset contains 8,000 news headlines categorized into four different topics (Rakib et al., 2020).
- **SearchSnippets:** Derived from web search



activities, it includes 12,340 search result snippets organized into eight distinct categories (Phan et al., 2008).

- **StackOverflow**: Comprising 20,000 question titles across 20 technical fields (Xu et al., 2017), this dataset is sampled from Kaggle competition data, covering technical discussions and programming-related queries.
- **Biomedical**: This dataset consists of 20,000 research paper titles in 20 scientific disciplines (Xu et al., 2017), sourced from BioASQ, showcasing the specialized terminology and format typical of academic research.
- **GoogleNews**: Providing a broad range of news content, it includes 11,109 articles related to 152 events (Yin and Wang, 2016). The dataset is available in three versions: complete articles (GoogleNews-TS), titles only (GoogleNews-T), and snippets only (GoogleNews-S).
- **Tweet**: Containing 2,472 tweets linked to 89 different queries (Yin and Wang, 2016), this dataset was gathered from the Text Retrieval Conference’s microblog tracks in 2011 and 2012, reflecting the casual and succinct nature of social media posts.

### B.3 Evaluation Metrics

Consistent with previous works (Rakib et al., 2020; Zheng et al., 2023), we employ two standard metrics to use the clustering performance: Accuracy (ACC) and Normalized Mutual Information (NMI). Accuracy measures the proportion of correct clustered texts, which is defined as:

$$ACC = \frac{\sum_{i=1}^N \mathbb{1}_{y_i = \text{map}(\tilde{y}_i)}}{N}, \quad (29)$$

where  $y_i$  is the true label and  $\tilde{y}_i$  is the predicted label,  $\text{map}(\cdot)$  operation refers to aligning the predicted labels with the true labels using the Hungarian algorithm. (Papadimitriou and Steiglitz, 1998).

Normalized Mutual Information (NMI) quantifies the shared information between the true and predicted label distributions, normalized by their individual uncertainties:

$$NMI(\mathbf{Y}, \tilde{\mathbf{Y}}) = \frac{I(\mathbf{Y}, \tilde{\mathbf{Y}})}{\sqrt{H(\mathbf{Y})H(\tilde{\mathbf{Y}})}} \quad (30)$$

where  $\mathbf{Y}$  and  $\tilde{\mathbf{Y}}$  represent the true and predicted label matrices respectively,  $I$  denotes mutual information, and  $H$  represents entropy.

### B.4 Experiment Settings

The batch size of the labeled and unlabeled data are set to  $B = 15$  and  $\mu B = 200$ , respectively. The temperature parameters for instance-level and prototypical-based contrastive learning are set to  $T_P = 1$  and  $T_I = 1$ . The outer loops of the GCG algorithm  $T_1$  and the iterations of the Lagrange multiplier algorithm  $T_2$  are set to 10. The total number of training iterations  $E_{total}$  is 1,500 for all datasets except the Tweet dataset, where  $E_{total} = 1,000$ . The number of first stage iterations  $E_{first}$  is 1,000 for all datasets except the Tweet dataset, in which  $E_{first} = 700$ . The maximum sentence length of the Encoder (BGE-M3) input is 32. The output dimension of the Projector  $h$  is set to  $D = 128$ .

## C Supplementary Experiment

### C.1 Clustering Degeneracy Study

We conducted comparative experiments to verify whether our method can prevent the occurrence of the clustering degeneracy problem. Clustering degeneracy is a significant challenge for imbalanced datasets (i.e., although the number of categories is provided to the model during training, the predicted number is still smaller than the real amount).

The results are shown in Figure 5. From these results, we can observe that, IOPC converges to the real category number, while other methods suffer from the clustering degeneracy problem.

### C.2 The visualization of text representations

To observe the distribution of samples in the feature space, we performed T-SNE visualization on SearchSnippets dataset for baseline models and IOPC. The result is shown in Figure 6. We can see that: (1) In **M3**, all the clusters overlap with each other. (2) **RSTC** shows some improvement over M3, but still contains a significant number of misclustered noise points, indicating poorer clustering performance. (3) **COTC** achieves a better representation distribution than RSTC, but it still has some errors, particularly confusing the clusters represented by red color and black color. (4) Our proposed **IOPC** achieves the best clustering performance. It effectively reduces the noise points within the clusters obtained by clustering. The

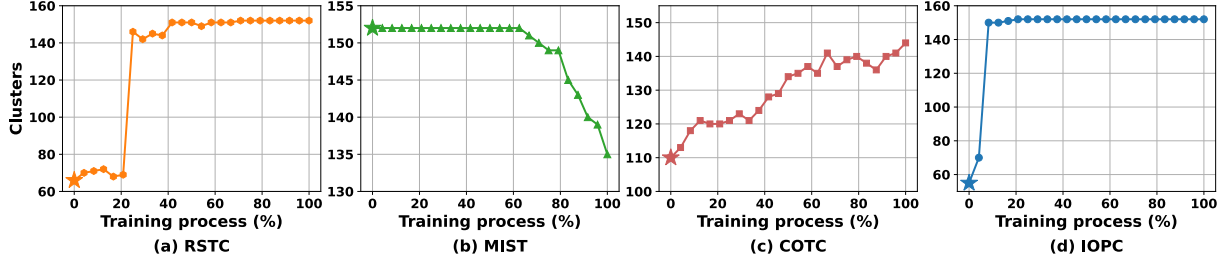


Figure 5: **Clustering Degeneracy Comparison.** The number of predicted clusters during the training process on the GoogleNews-T dataset.

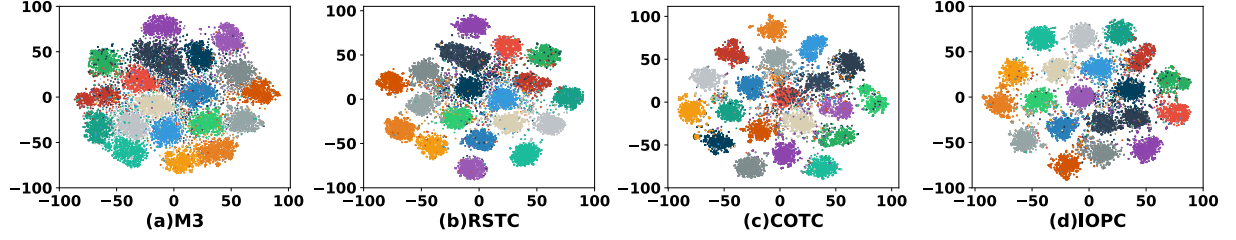


Figure 6: **T-SNE Comparison.** Each color indicates a truth category.

representation visualization indicates that our proposed method learned discriminative representations and achieved better clustering.

### C.3 The Comparison Results Using the Same Encoder

To ensure a fair comparison of algorithm performance, additional experiments were conducted using a unified Encoder. Among the baseline models, SCCL (Zhang et al., 2021), RSTC (Zheng et al., 2023), and COTC (Li et al., 2024) utilize the *distilbert-base-nli-stsb-mean-tokens* (SBERT) Encoder, MIST (Kamthawee et al., 2024) employs the *paraphrase-mpnet-base-v2* (MPNET) Encoder, and STSPL-SSC (Nie et al., 2024) uses the *bge-base-en-v1.5* (BGE-M3) Encoder. Notably, SBERT yields the lowest performance, MPNET surpasses SBERT, and BGE-M3 produces the best results.

In real-world short text clustering applications, the primary objective is to achieve the most accurate clustering results. To this end, IOPC adopts the same BGE-M3 Encoder used by STSPL-SSC (Nie et al., 2024). Different encoders may yield varying results; therefore, to ensure a fair comparison with previous studies, we replaced the encoders in the representative algorithms RSTC and COTC with the BGE-M3 Encoder.

The results, presented in Table 5, indicate that under identical Encoder conditions, IOPC continues to outperform the other models. Therefore, the superior performance achieved by IOPC is not closely related to the encoder.

### C.4 Research on Incorporating Labeled Data

Like the previous work STSPL-SSC, IOPC is a semi-supervised approach, while the other previous works are unsupervised methods. To ensure a fair comparison, we incorporated the same amount of labeled data used in IOPC into the previous works and applied the cross-entropy loss function to leverage the labeled data.

The results, presented in Table 6, indicate that simply incorporating a small amount of labeled data does not improve model performance. In fact, it has a negative impact. We attribute this to the fact that previous works utilize k-means to generate pseudo-labels at the beginning of the training process. K-means assigns random labels to the generated clusters, which may conflict with the true labels. Furthermore, these results demonstrate that the strong performance of our method is not solely due to the labeled data, but rather to its ability to effectively propagate knowledge from the labeled data to the unlabeled data.

### D Computation Budget

We built our model using PyTorch and performed all experiments on a NVIDIA GeForce RTX 3090 Ti GPU. To provide a comprehensive comparison with prior research, we evaluate both the parameter count and training time relative to existing methods, using the StackOverflow dataset as a benchmark. This comparison offers insights into the computational efficiency and scalability of our approach in relation to previous studies.

method	Agn	Sea	Sta	Bio	GN-TS	GN-T	GN-S	Twe
RSTC	89.39	81.26	86.78	51.67	84.21	80.12	82.82	77.06
STSPL-SSC	<u>89.92</u>	81.04	86.74	47.43	84.41	81.01	82.30	79.59
COTC	88.33	89.78	<u>89.83</u>	<u>51.92</u>	89.56	<u>85.02</u>	<u>87.10</u>	<u>91.53</u>
<b>IOPC</b>	<b>90.28</b>	<b>90.44</b>	<b>90.38</b>	<b>60.54</b>	<b>92.92</b>	<b>87.71</b>	<b>87.64</b>	<b>92.11</b>
improvement	<b>+0.36</b>	<b>+0.66</b>	<b>+0.55</b>	<b>+8.62</b>	<b>+3.36</b>	<b>+2.69</b>	<b>+0.54</b>	<b>+0.58</b>

Table 5: **Results of Using the Same Encoder** The experiment results for baseline models using the same Encoder.

method	Agn	Sea	Sta	Bio	GN-TS	GN-T	GN-S	Twe
RSTC	84.76	79.55	81.89	45.31	80.91	70.99	77.89	70.55
MIST	85.51	75.93	82.2	39.85	86.42	73.22	79.45	87.45
STSPL-SSC	<u>89.92</u>	81.04	86.74	47.43	84.41	81.01	82.30	79.59
COTC	87.06	<b>90.65</b>	<u>87.17</u>	<u>52.79</u>	<u>88.7</u>	<u>83.03</u>	<u>84.31</u>	<u>90.14</u>
<b>IOPC</b>	<b>90.28</b>	<u>90.44</u>	<b>90.38</b>	<b>60.54</b>	<b>92.92</b>	<b>87.71</b>	<b>87.64</b>	<b>92.11</b>
improvement	<b>+0.36</b>	<b>-0.21</b>	<b>+3.21</b>	<b>+7.75</b>	<b>+4.22</b>	<b>+4.68</b>	<b>+3.33</b>	<b>+1.97</b>

Table 6: **Results of Incorporating Labels for Baselines** The comparison between IOPC and previous models with labeled data incorporated.

	RSTC-origin	RSTC-M3	COTC-origin	COTC-M3	MIST-origin	<b>IOPC</b>
Training time	00:15:39	00:28:40	00:35:21	01:02:36	00:37:27	00:24:01
Parameters	68.25M	111.37M	77.44M	120.55M	109.5M	111.37M

Table 7: **The Comparison of Parameter Quantity and Training Time.** Where "RSTC-origin", "COTC-origin" and "MIST-origin" refer to the models presented in their respective original papers, while "RSTC-M3" and "COTC-M3" denote the models with the Encoder replaced by BGE-M3.

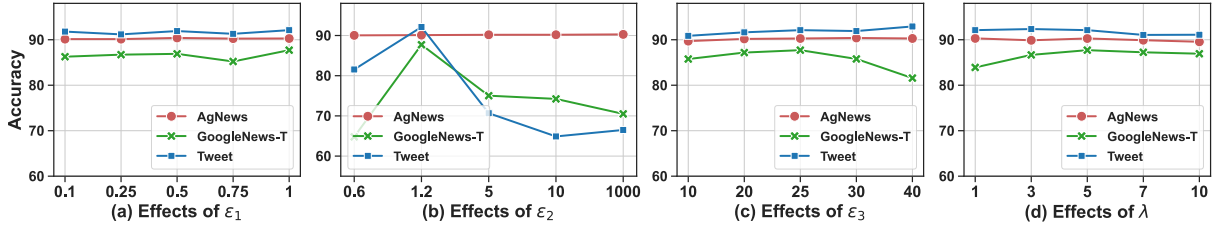


Figure 7: **Hyperparameter Analysis.** The effect of  $\epsilon_1$ ,  $\epsilon_2$ ,  $\epsilon_3$ , and  $\lambda$  on model accuracy.

The results in Table 7 show that, due to the adoption of BGE-M3 as the Encoder, our model has over more 40M parameters compared to RSTC-origin and COTC-origin. However, this increase is negligible relative to the significant improvement in clustering performance. Additionally, in previous work, MIST also uses a new Encoder, making its parameter count comparable to ours, but its clustering performance is still significantly lower than IOPC (as shown in Table 2). Furthermore, IOPC achieves the shortest training time except for RSTC-origin, indicating lower computational resource requirements. When RSTC and COTC are switched to BGE-M3 Encoder, their parameters and training time increase substantially.

## E Hyperparameter Analysis

We conducted a series of experiments to validate the effects of  $\epsilon_1$ ,  $\epsilon_2$ ,  $\epsilon_3$  and  $\lambda$  with values in  $\{0, 1, 5, 10\}$ ,  $\{0.03, 0.06, 0.1, 1, 3.5, 7, 10, 100\}$ ,  $\{10, 15, 20, 25, 30\}$  and  $\{1, 5, 10, 15, 20\}$ , respectively. The experiments were conducted on the representative datasets **AgNews**, **GoogleNews-T** and **Tweet**, with results presented in Figure 7.

From Figures 7(a), 7(c), and 7(d), we observe that variations in  $\epsilon_1$ ,  $\epsilon_3$ , and  $\lambda$  have minimal impact on model performance, suggesting that the model is largely insensitive to these parameters. In contrast, Figure 7(b) emphasizes the importance of tuning  $\epsilon_2$  for imbalanced datasets, whereas it has no discernible effect on balanced datasets. Since

$\varepsilon_2$  regulates the penalty strength for the imbalance levels of predicted cluster probabilities in Eq.(10), we determine its value based on the degree of imbalance in the dataset.

Although our model has several hyperparameters, only  $\varepsilon_2$  influences the performance on imbalanced datasets. This suggests that the model exhibits strong robustness and generalizability, as it remains largely unaffected by other hyperparameter variations. Consequently, when applied to unseen data, the model demonstrates higher adaptability, requiring minimal hyperparameter tuning for effective performance. Experimentally, we set  $\varepsilon_1 = 1$ ,  $\varepsilon_3 = 25$  and  $\lambda = 5$  for all datasets;  $\varepsilon_2 = 1000$  and 1.2 for balanced datasets and severely imbalanced datasets, respectively.