

Label-enhanced Prototypical Network with Contrastive Learning for Multi-label Few-shot Aspect Category Detection

Han Liu
Dalian University of Technology
Dalian, China
liu.han.dut@gmail.com

Feng Zhang
Peking University
Beijing, China
zhangfeng@stu.pku.edu.cn

Xiaotong Zhang*
Dalian University of Technology
Dalian, China
zxt.dut@hotmail.com

Siyang Zhao
Dalian University of Technology
Dalian, China
zhao_siyang@mail.dlut.edu.cn

Junjie Sun
Dalian University of Technology
Dalian, China
sunjunjiedlut@hotmail.com

Hong Yu
Dalian University of Technology
Dalian, China
hongyu@dlut.edu.cn

Xianchao Zhang*
Dalian University of Technology
Dalian, China
xczhang@dlut.edu.cn

ABSTRACT

Multi-label aspect category detection allows a given review sentence to contain multiple aspect categories, which is shown to be more practical in sentiment analysis and attracting increasing attention. As annotating large amounts of data is time-consuming and labor-intensive, data scarcity occurs frequently in real-world scenarios, which motivates multi-label few-shot aspect category detection. However, research on this problem is still in infancy and few methods are available. In this paper, we propose a novel label-enhanced prototypical network (LPN) for multi-label few-shot aspect category detection. The highlights of LPN can be summarized as follows. First, it leverages label description as auxiliary knowledge to learn more discriminative prototypes, which can retain aspect-relevant information while eliminating the harmful effect caused by irrelevant aspects. Second, it integrates with contrastive learning, which encourages that the sentences with the same aspect label are pulled together in embedding space while simultaneously pushing apart the sentences with different aspect labels. In addition, it introduces an adaptive multi-label inference module to predict the aspect count in the sentence, which is simple yet effective. Extensive experimental results on three datasets demonstrate that our proposed model LPN can consistently achieve state-of-the-art performance.

CCS CONCEPTS

• **Computing methodologies** → **Artificial intelligence; Natural language processing.**

KEYWORDS

Multi-label Few-shot Learning; Aspect Category Detection; Prototypical Network; Contrastive Learning

ACM Reference Format:

Han Liu, Feng Zhang, Xiaotong Zhang, Siyang Zhao, Junjie Sun, Hong Yu, and Xianchao Zhang. 2022. Label-enhanced Prototypical Network with Contrastive Learning for Multi-label Few-shot Aspect Category Detection. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22)*, August 14–18, 2022, Washington, DC, USA. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3534678.3539340>

1 INTRODUCTION

Aspect Category Detection (ACD) is a fundamental task in sentiment analysis, which aims to identify the aspect categories mentioned in a given review sentence from a predefined aspect category set. As human usually make comments from different angles, i.e., a review sentence always contains multiple aspects, multi-label aspect category detection task came into existence. Existing approaches for multi-label ACD have achieved impressive and promising performance [17, 20]. However, they rely heavily on large amounts of labeled data for each aspect. As annotating data is usually time-consuming, labor-intensive and even unachievable in real-world application, which motivates the multi-label few-shot aspect category detection task.

Few-shot learning can recognize novel categories effectively with only a handful of labeled samples by exploiting the prior knowledge learned from previous categories, which is promising to break the data-shackles. Recent methods have made great progress in computer vision domain [22, 37] and natural language processing domain [10, 14]. Among these methods, prototypical network [30] is a powerful and potential model, which follows the episode learning strategy and uses the N -way K -shot setting. Specifically, in each

*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD '22, August 14–18, 2022, Washington, DC, USA

© 2022 Association for Computing Machinery.
ACM ISBN 978-1-4503-9385-0/22/08...\$15.00
<https://doi.org/10.1145/3534678.3539340>

episode, prototypical network first learns each prototype by averaging the corresponding K support samples, and then predicts the labels of query samples based on the negative Euclidean distance between query samples and N prototypes in support set.

Intuitively, we can directly extend prototypical network to solve the multi-label few-shot aspect category detection problem. However, there exist several challenging issues. (1) Simply calculating the prototype by averaging intra-class support samples may cause that different aspects share an identical prototype. Take an example in Figure 1, we construct an episode under 3-way 2-shot setting, i.e., there are 3 aspects with 2 samples per aspect in support set. As a review sentence may contain multiple aspects, it is quite possible that different aspects distribute in the same support samples. Here "staff" and "food" have the same support samples. Thus we will obtain the same prototype for "staff" and "food". Although the work [16] attempts to use the attention mechanism to alleviate this issue, we find that it still does not work in this case. (2) When learning the prototype for a target aspect in multi-label scenarios, as each sentence probably contains several aspects, some irrelevant aspects will inevitably disturb the learning procedure. For example in Figure 1, when learning the prototype for "staff", the irrelevant aspects like "food" and "location" in "The views are amazing from any location, staff is friendly and the food was great too!" will cause some negative impact. (3) Due to the diversity of human expression, different sentences may contain different numbers of aspects, so it is urgently needed to design an effective model to automatically predict the number of aspects in a sentence.

In this paper, we propose a novel label-enhanced prototypical network (LPN) for multi-label few-shot aspect category detection. The main contributions of LPN consist of three parts. (1) By utilizing the label text description as complementary information to calculate the relationship between sentences and aspect labels and then obtain more discriminative prototypes, the LPN model can not only avoid the issue that different aspects share an identical prototype, but also retain aspect-relevant information while eliminating the negative effect triggered by irrelevant (noisy) aspects. (2) By integrating contrastive learning to obtain more powerful embeddings, the LPN model can push the embeddings from the same class close and embeddings from different classes further apart, thus facilitating the downstream aspect category detection task. (3) By introducing the adaptive multi-label inference module, the LPN model can determinate the number of aspects accurately. To verify the effectiveness of our proposed model, we conduct extensive experiments on three datasets. The empirical study shows that LPN can achieve state-of-the-art performance in comparison with other strong baselines.

2 RELATED WORK

2.1 Aspect Category Detection

Aspect Based Sentiment Analysis (ABSA) [33] is a fine-grained sentiment analysis task that aims to extract aspects and predict the sentiment of each aspect. Aspect category detection (ACD) is an important subtask of ABSA, which aims to categorize a given review sentence into a set of predefined aspects. Previous studies mainly focus on single-aspect category detection, which include unsupervised and supervised methods. Unsupervised methods use

Support Set	
staff	(1) It is the staff and food quality that really needs fixing. (2) The views are amazing from any location , staff is friendly and the food was great too!
food	(1) It is the staff and food quality that really needs fixing. (2) The views are amazing from any location , staff is friendly and the food was great too!
experience	(1) Incredible spa experience! (2) The food is always good and service has always been a great experience .
Query Set	
experience and staff	(1) It was such a horrible experience , she was rude, unmannered and non professional great clip should not retain such a waste employee!
staff and food	(2) The pool is gorgeous, the rooms clean, delicious food , and staff that went above and beyond to help us enjoy our stay.
food	(3) We had breakfast the next morning on the first floor and the food was surprisingly good.

Figure 1: A meta-task example in 3-way 2-shot setting. The first column denotes the aspect label and the second column denotes the corresponding review sentence. As each review sentence may contain multiple aspects, we use different color background to mark the key words. The words in gray describe irrelevant (noisy) aspects, and the words in other colors represent the target aspects.

semantic association analysis based on point-wise mutual information [31] or co-occurrence association rule mining [29] to extract aspects. These methods require a large amount of corpus resources and the performance is also barely satisfactory. Supervised methods exploit representation learning [42], topic-attention network [24] or multilingual ngram-based convolutional network [12] to identify different aspect categories. These methods have shown promising results in practice, but they rely heavily on a considerable amount of labeled data for each aspect to train a discriminative classifier. Due to the diversity and casualness of human expression, a review sentence often contains multiple aspects, which motivates multi-aspect category detection. Existing approaches for multi-label ACD [17, 20] have achieved impressive performance. However, similar with supervised methods for single-aspect category detection, they also suffer from the serious data scarcity issue.

2.2 Few-shot Learning

Few-shot learning is a paradigm for solving the data deficiency problem, which aims to use the knowledge learned from seen classes, of which abundant labeled samples are available for training, to recognize unseen classes, of which limited labeled samples are provided. It has drawn much attention in computer vision domain [22, 37] and natural language processing domain [18, 40]. Meta-learning has been successfully applied to solve the few-shot learning problem, which mainly includes model-based approaches, optimization-based approaches and metric-based approaches. Specifically, for model-based methods, like MANN [27] and MetaNet [25], they depend on the models which can update the parameters rapidly with a few training steps. For optimization-based methods, like LSTM

Meta-Learner [26], MAML [8], Bayesian MAML [39], they intend to adjust some optimization algorithms so that the model can be good at learning with a few examples. For metric-based methods, such as matching network [35], prototypical network [30], relation network [32] and so on [3, 11], their basic idea is to learn a feature mapping function that projects support and query samples into an embedding space and classify the queries by learning their relations by some metrics in that space. Among these methods, due to the simplicity and effectiveness, prototypical network is one of the most popular methods in few-shot learning.

2.3 Multi-label Few-shot Learning

Traditional few-shot learning focuses on single-label classification task. However, in many real scenarios a sample often has multiple labels, which gives birth to multi-label few-shot learning. As far as we know, only a few works have been done for this task. In computer vision domain, LaSO [1] is a multi-label few-shot image classification model which leverages the label set operations (intersection, union, subtraction) to guide the model to learn the semantic features. In audio domain, Cheng et al. [6] use the one-versus-rest episode selection strategy and attention mechanism to deal with the multi-label few-shot sound event recognition problem. In natural language processing domain, Hou et al. [15] focus on multi-label few-shot intent classification task and propose a meta calibrated threshold mechanism with kernel regression and logits adapting that estimates threshold using both prior domain experience and new domain knowledge.

Proto-AWATT [16] is the first work which aims to address aspect category detection in the few-shot scenario. It attempts to leverage support-set and query-set attention mechanisms to alleviate the negative effect caused by noisy aspects, and has achieved the state-of-the-art performance. However, it still suffers from the issue that different aspects perhaps share an identical prototype. In addition, Proto-AWATT learns a dynamic threshold via the policy network, which requires a more complex two-stage training process and that the threshold satisfies the idealized Beta distribution assumption.

3 PROBLEM FORMULATION

To ease understanding, we briefly introduce the task of multi-label few-shot aspect category detection. Table 1 summarizes some symbol explanation in details.

Few-shot learning aims to recognize unknown categories with few labeled samples by leveraging prior knowledge learned from known categories. In general, the data can be divided into two parts: seen (known) class set C_{seen} and unseen (unknown) class set C_{unseen} , and $C_{\text{seen}} \cap C_{\text{unseen}} = \emptyset$. A classifier is trained with numerous samples from C_{seen} , and it is quickly adopted to C_{unseen} (which is unavailable in training) with only a few labeled data. Meta learning is an effective solution for few-shot learning, which contains two phases: meta-training and meta-testing, and it commonly follows the N -way K -shot setting, i.e., for each task, there are N classes and each class has K supports (labeled samples).

In meta-training phase, the meta-classifier is trained on N_{train} tasks. In each training task, it consists of a support set and a query set. To construct the training task, N classes are randomly sampled from N_{seen} . The support set is composed of randomly selecting K

Table 1: Symbol explanation.

Symbol	Explanation
C_{seen}	the seen class set
C_{unseen}	the unseen class set
N	the number of aspect classes in each episode
K	the number of support shots in each class
S	the support set of an episode
Q	the query set of an episode
\mathbf{x}	a sentence with T words
y	the class label of \mathbf{x}
H	the embedding matrix of \mathbf{x} via any pre-trained model
\mathbf{o}	the representation of \mathbf{x} via feature extraction
E	the label description representation
\mathbf{p}^i	the label-enhanced prototype of class i
\mathbf{z}^i	the label-specific embedding associated with class i

labeled samples from each of the N classes, i.e., $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$, where \mathbf{x}_i is a data sample, \mathbf{y}_i is the class label and $m = N \times K$. The query set consists of a portion of the remaining samples from these N classes, i.e., $Q = \{(\mathbf{x}_j, \mathbf{y}_j)\}_{j=1}^n$, where n is the number of queries.

In meta-testing phase, the trained meta-classifier is used to predict the labels of queries in N_{test} tasks. In each test task, it also has a support set and a query set. In a similar manner, N classes are randomly sampled from the test classes C_{unseen} . The support set and query set are constructed in the same way as those in meta-training phase. As the labels of queries are unknown in testing stage, the query set in the test task can be represented as $Q = \{\mathbf{x}_j\}_{j=1}^n$. The goal is to predict the class labels for these queries.

Multi-label few-shot aspect category detection allows that each single sentence is associated with a set of aspect categories simultaneously. Specifically, given a sentence \mathbf{x} , its label can be represented with a vector $\mathbf{y} = \{y^1, y^2, \dots, y^N\} \in \mathbb{R}^N$, where $y^i \in \{0, 1\}$ and N is the number of possible aspects. In this paper, we focus on the multi-label few-shot aspect category detection problem.

4 APPROACH

The overall framework of the proposed LPN is shown in Figure 2. It consists of four components: feature extraction, label-enhanced prototypical network, contrastive learning and adaptive multi-label inference. In this section, we will introduce these modules in details.

4.1 Feature Extraction

Given a sentence \mathbf{x} with T words, we can use any pre-trained language model like Bert [7] to encode each word (token), and then get the embedding matrix $H = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_T] \in \mathbb{R}^{d \times T}$. To better extract the sentence-level semantic feature and assign reasonable importance for each word, we follow [21, 38] to utilize a multi-head self-attentive module to generate the sentence embedding. Specifically,

$$\mathbf{A} = \text{softmax}(F_2 \tanh(F_1 H)), \quad (1)$$

where $\mathbf{A} \in \mathbb{R}^{R \times T}$ is the self-attention weight matrix, R is the number of independent attention heads, $F_1 \in \mathbb{R}^{d' \times d}$ and $F_2 \in \mathbb{R}^{R \times d'}$ are trainable parameter matrices. After obtaining \mathbf{A} , we first calculate

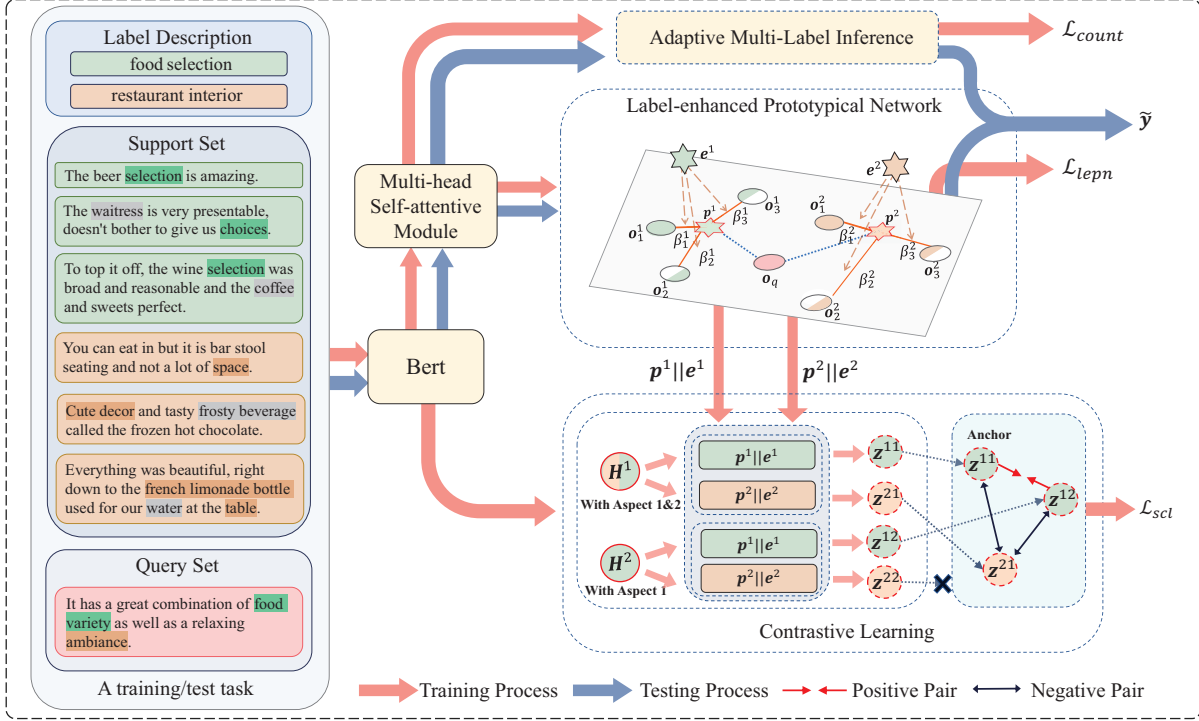


Figure 2: Illustration of our proposed method LPN.

the embedding matrix by:

$$M = HA^T, \quad (2)$$

where $M = [m_1, m_2, \dots, m_R] \in \mathbb{R}^{d \times R}$. Furthermore, we concatenate the obtained embeddings from different heads, and use a simple linear projection to calculate the embedding of the sentence,

$$o = F_3[m_1 || m_2 || \dots || m_R], \quad (3)$$

where $F_3 \in \mathbb{R}^{d \times dR}$ is a trainable parameter matrix, and $||$ represents the concatenation operation. $o \in \mathbb{R}^d$ is the final representation of the sentence x .

4.2 Label-enhanced Prototypical Network

In terms of multi-label few-shot aspect category detection task, the key point is to learn more discriminative prototypes which could better retain class-relevant information while eliminating the harmful effect caused by other noisy (irrelevant) aspects. To achieve this goal, we propose to leverage the label text description information to calculate the relationship between sentences and aspect labels, thus obtaining more representative prototypes.

Considering the N -way K -shot setting, we have a support set S , which can be represented by $S = \{x_1^1, x_2^1, \dots, x_K^1, \dots, x_1^N, x_2^N, \dots, x_K^N\}$, where x_j^i denotes the j -th sample belonging to the i -th class. After feature extraction, we get the representations of these samples $O = \{o_1^1, o_2^1, \dots, o_K^1, \dots, o_1^N, o_2^N, \dots, o_K^N\}$. In a similar manner, for each label description like "Room cleanliness" or "Staff owner", we can obtain its corresponding representation via feature extraction. In

the N -way scenario, the label description representation can be represented as $E = \{e^1, e^2, \dots, e^N\}$.

When calculating the class prototype, as each sentence may contain multiple aspects, we would better first determinate the importance weight of each sentence for a class prototype. To this end, we utilize the label text description as auxiliary information. Specifically,

$$\alpha_j^i = \sigma_j^{i,T} W e^i, \quad (4)$$

where α_j^i denotes the importance weight of the j -th sentence for the i -th class prototype. $\sigma_j^i \in \mathbb{R}^d$ denotes the representation of the j -th sample belonging to the i -th class, $\sigma_j^{i,T} \in \mathbb{R}^{1 \times d}$ is the transpose of σ_j^i . $W \in \mathbb{R}^{d \times d}$ is a trainable projection matrix. e^i is the representation of the i -th label description.

Inspired by low-rank bilinear model [41], if imposing a low-rank restriction on W , Eq. (4) can be rewritten as follows:

$$\alpha_j^i = \sigma_j^{i,T} UV^T e^i = \mathbf{1}^T (U^T \sigma_j^i \circ V^T e^i), \quad (5)$$

where $U \in \mathbb{R}^{d \times k}$ and $V \in \mathbb{R}^{d \times k}$ are two low-rank matrices with $k < d$. $\mathbf{1}^T$ is a all-one vector. \circ is the Hadamard product, i.e., element-wise multiplication. By using this low-rank trick, we can reduce the number of parameters to some extent.

To make the coefficients comparable among different sentences, we normalize them across K sentences (shots) belonging to the same class with the softmax function:

$$\beta_j^i = \frac{\exp(\alpha_j^i)}{\sum_{j'=1}^K \exp(\alpha_{j'}^i)}. \quad (6)$$

Then we calculate the label-enhanced prototype $\mathbf{p}^i \in \mathbb{R}^d$ for class i by:

$$\mathbf{p}^i = \sum_{j=1}^K \beta_j^i \mathbf{o}_j. \quad (7)$$

Given a query sentence $\mathbf{x} \in \mathcal{Q}$, we can compute the conditional probability $p(y = y^i | \mathbf{x}, \mathcal{S})$ to predict its aspect label based on negative squared Euclidean distance.

$$p(y = y^i | \mathbf{x}, \mathcal{S}) = \frac{\exp(-\|\mathbf{o} - \mathbf{p}^i\|_2^2)}{\sum_{j=1}^N \exp(-\|\mathbf{o} - \mathbf{p}^j\|_2^2)}, \quad (8)$$

where \mathbf{o} denotes the representation of \mathbf{x} , which is obtained via feature extraction.

Finally, we perform the cross-entropy loss on all samples in the query set \mathcal{Q} , i.e., the loss function of label-enhanced prototypical network \mathcal{L}_{lep} can be written as:

$$\mathcal{L}_{lep} = \frac{1}{|\mathcal{Q}|} \sum_{\mathbf{x} \in \mathcal{Q}} \sum_{i=1}^N -y^i \log p(y = y^i | \mathbf{x}, \mathcal{S}), \quad (9)$$

where $|\mathcal{Q}|$ is the number of samples in \mathcal{Q} . Note that in the multi-label N -way K -shot setting, as a sentence may have multiple labels, we need to consider N labels for each query sentence.

4.3 Integrating with Contrastive Learning

Contrastive learning has achieved great success in computer vision [5, 36], which aims to maximize similarities between instances from the same class and minimize similarities between instances from different classes. Here we integrate the contrastive learning into label-enhanced prototypical network to generate better sentence embeddings.

For traditional single-label aspect detection, we can directly construct the contrastive samples using the known aspect labels. However, in the multi-label aspect detection scenario, as a sentence may contain a couple of aspects, for example, one sentence "The pool is gorgeous, the room clean, delicious food, and staff that went above and beyond to help us enjoy our stay" contains two aspect labels "food" and "staff", and the other sentence "The food is always good and service has always been a great experience" contains two aspect labels "food" and "experience". If simply treating these two sentences as positive pairs, it is unreasonable obviously. The reason is that though these two sentences share a common aspect label "food", they also have a totally different aspect label. To alleviate the above issue, we use the prototypes and label description information to first generate the label-specific embeddings for each sentence, and then construct the contrastive samples.

In the N -way K -shot setting, for each meta-task, we can obtain N prototypes $\mathbf{P} = \{\mathbf{p}^1, \mathbf{p}^2, \dots, \mathbf{p}^N\}$ and N label description representations $\mathbf{E} = \{\mathbf{e}^1, \mathbf{e}^2, \dots, \mathbf{e}^N\}$. By combining \mathbf{P} and \mathbf{E} , we can have the prototypes integrated with label description information $\{\mathbf{a}^1, \mathbf{a}^2, \dots, \mathbf{a}^N\}$, where $\mathbf{a}^i = [\mathbf{p}^i || \mathbf{e}^i] \in \mathbb{R}^{2d}$ and $||$ represents the concatenation operation. Then for a sentence \mathbf{x} in the meta-task, we can compute its label-specific embedding $\mathbf{z}^i \in \mathbb{R}^d$ associated with label i by:

$$\mathbf{z}^i = \mathbf{g}^i \mathbf{H}^T, \quad (10)$$

where $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_T] \in \mathbb{R}^{d \times T}$ is the embedding matrix obtain from any pre-trained language model like Bert. $\mathbf{g}^i \in \mathbb{R}^{1 \times T}$ is a

weight vector obtained by:

$$\mathbf{g}^i = \text{softmax}((\mathbf{W}_a \mathbf{a}^i + \mathbf{b}_a)^T \mathbf{H}), \quad (11)$$

where $\mathbf{W}_a \in \mathbb{R}^{d \times 2d}$ and $\mathbf{b}_a \in \mathbb{R}^d$ are trainable parameters.

In meta-training phase, given a meta-task with N_t samples, we first use Eq. (10) and (11) to get the label-specific embeddings and then collect all these embeddings to construct the set $\mathbf{Z} = \{\mathbf{z}^{ij} \in \mathbb{R}^d | i \in \{1, 2, \dots, N\}, j \in \{1, 2, \dots, N_t\}\}$. If we regard each label-specific embedding as an independent instance, each \mathbf{z}^{ij} will be associated with a single ground-truth label y^{ij} . Specifically, a sentence "The food is always good and service has always been a great experience" contains two aspect labels "food" and "experience", the labels of obtained label-specific embedding associated with "food" and "experience" will be set to 1. Then we can define the set $\mathbf{Y} = \{y^{ij} \in \{0, 1\} | i \in \{1, 2, \dots, N\}, j \in \{1, 2, \dots, N_t\}\}$. Furthermore, we define the set $\mathbf{I} = \{\mathbf{z}^{ij} \in \mathbf{Z} | y^{ij} = 1\}$ which contains the label-specific embeddings with usable ground-truth labels, and the set $\mathbf{\Gamma}^{ij} = \{\mathbf{I} \setminus \mathbf{z}^{ij}\}$ which contains the embeddings in \mathbf{I} with \mathbf{z}^{ij} excluded. Considering \mathbf{z}^{ij} as the anchor, we can generate the positive sample set $\mathbf{\Lambda}^{ij} = \{\mathbf{z}^{ik} \in \mathbf{\Gamma}^{ij} | y^{ik} = y^{ij} = 1\}$ for \mathbf{z}^{ij} , and the negative samples for \mathbf{z}^{ij} are the remaining ones in $\mathbf{\Gamma}^{ij}$. With above notations, the contrastive learning loss for the anchor \mathbf{z}^{ij} can be written as:

$$\mathcal{L}_{scl}^{ij} = -\frac{1}{|\mathbf{\Lambda}^{ij}|} \sum_{\mathbf{z}^{ik} \in \mathbf{\Lambda}^{ij}} \log \frac{\exp(\mathbf{z}^{ij} \cdot \mathbf{z}^{ik} / \tau)}{\sum_{\mathbf{z}^* \in \mathbf{\Gamma}^{ij}} \exp(\mathbf{z}^{ij} \cdot \mathbf{z}^* / \tau)}, \quad (12)$$

where $\mathbf{z}^{ij} \cdot \mathbf{z}^{ik}$ denotes the inner product of the two vectors. $|\mathbf{\Lambda}^{ij}|$ is the number of embeddings in $\mathbf{\Lambda}^{ij}$. $\tau > 0$ is an adjustable scalar parameter, which can control the separation degree of classes [13]. By considering all the anchors, we can have the entire contrastive loss as follows.

$$\mathcal{L}_{scl} = \frac{1}{|\mathbf{I}|} \sum_{\mathbf{z}^{ij} \in \mathbf{I}} \mathcal{L}_{scl}^{ij}. \quad (13)$$

To analyze Eq. (13), we do some simple formula transformation as below.

$$\begin{aligned} \mathcal{L}_{scl} &= \frac{1}{|\mathbf{I}|} \sum_{\mathbf{z}^{ij} \in \mathbf{I}} -\frac{1}{|\mathbf{\Lambda}^{ij}|} \mathcal{L}', \\ \mathcal{L}' &= \sum_{\mathbf{z}^{ik} \in \mathbf{\Lambda}^{ij}} \log \frac{\exp(\mathbf{z}^{ij} \cdot \mathbf{z}^{ik} / \tau)}{\sum_{\mathbf{z}^* \in \mathbf{\Gamma}^{ij}} \exp(\mathbf{z}^{ij} \cdot \mathbf{z}^* / \tau)} \\ &= \sum_{\mathbf{z}^{ik} \in \mathbf{\Lambda}^{ij}} \underbrace{\left(\frac{\mathbf{z}^{ij} \cdot \mathbf{z}^{ik}}{\tau} \right)}_{\text{positive}} - \log \underbrace{\left(\sum_{\mathbf{z}^* \in \mathbf{\Gamma}^{ij}} \exp\left(\frac{\mathbf{z}^{ij} \cdot \mathbf{z}^*}{\tau}\right) \right)}_{\text{positive+negative}}. \end{aligned} \quad (14)$$

From the above formula, it is easy to find that if we want to minimize \mathcal{L}_{scl} , we must maximize \mathcal{L}' , where we need to maximize the positive term and minimize the positive+negative term, so the negative term will be decreased. Intuitively, the contrastive learning technique can push the label-specific embeddings from the same class close and embeddings from different classes further apart.

4.4 Adaptive Multi-label Inference

For multi-label few-shot aspect category detection, one of the challenges is to determine the number of aspects in the sentence. Previous work [16] learns a dynamic threshold via the policy network.

Table 2: Dataset statistics. #Aspects and #Sentences denote the number of aspects and sentences respectively.

Dataset	Split	#Aspects	#Sentences
FewAsp (single)	Training	64	12800
	Validation	16	3200
	Testing	20	4000
FewAsp (multi)	Training	64	25600
	Validation	16	6400
	Testing	20	8000
FewAsp	Training	64	40320
	Validation	16	10080
	Testing	20	12600

However, it requires that the threshold satisfies the Beta distribution assumption, which seems a little over-idealized. In addition, as it is a two-stage method, the training process is also more complicated. To overcome this issue, we propose an adaptive multi-label inference method, which is simple yet effective.

In the N -way K -shot setting, N is the maximal number of aspects in a sentence. Given a sentence \mathbf{x} , we can get its representation $\mathbf{o} \in \mathbb{R}^d$ via feature extraction. Then we use a multi-layer perception to predict the number of aspects in \mathbf{x} . Specifically,

$$\mathbf{n}_l = \text{softmax}(\mathbf{W}_l \mathbf{o} + \mathbf{b}_l), \quad (15)$$

where $\mathbf{W}_l \in \mathbb{R}^{N \times d}$ and $\mathbf{b}_l \in \mathbb{R}^N$ are trainable parameters. $\mathbf{n}_l \in \mathbb{R}^N$ is the indicator for the number of aspects. Take an example, if the maximal value of \mathbf{n}_l is the second element, it means that \mathbf{x} contains two aspects.

Then in the meta-training stage, for each sentence from support set \mathcal{S} and query set \mathcal{Q} , we use cross entropy to calculate the loss of the aspect count,

$$\mathcal{L}_{count} = \frac{1}{|\mathcal{S} \cup \mathcal{Q}|} \sum_{\mathbf{x} \in \mathcal{S} \cup \mathcal{Q}} -\mathbf{1}^T (\mathbf{t}_l \circ \log(\mathbf{n}_l)), \quad (16)$$

where $\mathbf{1}^T$ is a all-one vector. \circ is the Hadamard product, i.e., element-wise multiplication. $\log(\mathbf{n}_l) \in \mathbb{R}^N$ is to do the log operation on each element of \mathbf{n}_l . \mathbf{t}_l is the ground-truth aspect count vector of \mathbf{x} .

By combining Eq. (9), (13) and (16), we have the overall loss function of the proposed framework:

$$\mathcal{L}_{total} = \mathcal{L}_{lepn} + \gamma \mathcal{L}_{scl} + \lambda \mathcal{L}_{count}, \quad (17)$$

where γ and λ are adjustable trade-off parameters. By minimizing \mathcal{L}_{total} with the gradient descent method, all trainable parameters can be learned.

5 EXPERIMENTS

5.1 Datasets

For fair comparison, we exactly follow [16] to perform experiments on three datasets: FewAsp (single), FewAsp (multi) and FewAsp. All these datasets are sampled from the large-scale multi-domain dataset for aspect recommendation YelpAspect [4]. Specifically, FewAsp (single) consists of single-aspect sentences, FewAsp (multi) consists of a majority of multi-aspect sentences with a minority of

Table 3: Hyperparameters of our proposed method LPN.

Model	d	d'	R	k	λ	γ	τ
LPN	768	256	4	100	0.1	0.01	0.1

single-aspect sentences, as some aspects only have a small amount of multi-aspect samples, and FewAsp is randomly sampled from the original dataset, which follows the same data distribution with the real scenario. For data split, we also follow [16] to divide the 100 aspects without intersection into 64 aspects for training, 16 aspects for validation, and 20 aspects for testing. The detailed dataset statistics is shown in Table 2.

5.2 Baselines

We compare the proposed LPN model with the following strong baselines: Matching Network [35], Relation Network [32], Graph Network [28], Prototypical Network [30], IMP [2], Proto-HATT [9] and Proto-AWATT [16].

- **Matching Network** [35] first learns a embedding mapping function and then takes the cosine similarity as distance measure to obtain the classification results.
- **Prototypical Network** [30] calculates the prototype for each class by averaging the corresponding support samples, and utilizes the negative Euclidean distance between query samples and prototypes to do the few-shot classification task.
- **Relation Network** [32] uses a deep neural network instead of the fixed distance measure to calculate the relationship between query and support samples.
- **Graph Network** [28] attempts to cast few-shot learning as a supervised message passing task which is trained end-to-end using graph neural networks.
- **IMP** [2] introduces infinite mixture prototypes for few-shot learning, which represents each class by a set of clusters.
- **Proto-HATT** [9] is a hybrid attention-based prototypical networks for the problem of noisy few-shot relation classification. It uses instance-level and feature-level attention schemes to highlight the crucial instances and features respectively.
- **Proto-AWATT** [16] is the first method for multi-label few-shot aspect category detection task. It utilizes support-set and query-set attention mechanisms to alleviate the adverse effect caused by noisy aspects.

In addition, we conduct ablation study to evaluate the contribution of label enhancement and contrastive learning in LPN. Specifically, we evaluate LPN under three cases: LPN (o, o), LPN (w, o) and LPN (w, w).

- **LPN (o, o)** means the LPN model without label enhancement and contrastive learning.
- **LPN (w, o)** means the LPN model with label enhancement, but without contrastive learning.
- **LPN (w, w)** means the LPN model with label enhancement and contrastive learning, i.e., the final model.

Table 4: Average AUC and macro-F1 score on FewAsp (single).

Model	5-way 5-shot		5-way 10-shot		10-way 5-shot		10-way 10-shot	
	AUC	F1	AUC	F1	AUC	F1	AUC	F1
Matching Network [35]	97.05	81.89	97.49	84.62	96.30	70.95	96.72	73.28
Prototypical Network [30]	96.49	83.30	97.53	86.29	95.97	74.23	96.71	76.83
Relation Network [32]	93.31	75.79	90.86	72.02	91.81	63.78	90.54	61.15
Graph Network [28]	96.54	81.45	97.46	85.04	95.45	70.75	96.97	77.84
IMP [2]	96.65	83.69	97.47	86.14	96.00	73.80	96.91	77.09
Proto-HATT [9]	96.45	83.33	97.62	86.71	95.71	73.42	97.00	77.65
Proto-AWATT [16]	97.56	86.71	97.96	88.54	97.01	80.28	97.55	82.97
LPN (o, o)	97.88	87.62	98.48	90.31	98.13	83.99	98.53	85.95
LPN (w, o)	99.22	92.61	99.35	93.57	99.11	89.35	99.32	91.08
LPN (w, w)	99.29	94.43	99.49	94.40	99.14	89.40	99.28	90.43

Table 5: Average AUC and macro-F1 score on FewAsp (multi).

Model	5-way 5-shot		5-way 10-shot		10-way 5-shot		10-way 10-shot	
	AUC	F1	AUC	F1	AUC	F1	AUC	F1
Matching Network [35]	89.54	65.70	91.38	69.02	88.28	50.86	89.94	54.42
Prototypical Network [30]	89.67	67.88	91.60	72.32	88.01	52.72	90.68	58.92
Relation Network [32]	84.91	58.38	86.21	61.37	84.22	43.71	84.72	44.85
Graph Network [28]	87.97	59.25	90.45	64.63	86.05	45.42	88.44	48.49
IMP [2]	90.12	68.86	92.29	73.51	88.71	53.96	91.10	59.86
Proto-HATT [9]	91.10	69.15	93.03	73.91	90.44	55.34	92.38	60.21
Proto-AWATT [16]	91.45	71.72	93.89	77.19	89.80	58.89	92.34	66.76
LPN (o, o)	93.09	72.45	94.92	76.89	92.95	61.33	94.62	66.39
LPN (w, o)	95.43	78.82	96.22	81.70	94.29	66.36	95.43	71.08
LPN (w, w)	95.66	79.48	96.55	82.81	94.51	67.28	95.66	71.87

Table 6: Average AUC and macro-F1 score on FewAsp.

Model	5-way 5-shot		5-way 10-shot		10-way 5-shot		10-way 10-shot	
	AUC	F1	AUC	F1	AUC	F1	AUC	F1
Matching Network [35]	90.76	67.14	92.39	70.09	88.44	51.27	89.90	54.61
Prototypical Network [30]	88.88	66.96	91.77	73.27	87.35	52.06	90.13	59.03
Relation Network [32]	85.56	59.52	86.98	62.78	84.94	45.62	83.77	44.70
Graph Network [28]	89.48	61.49	92.35	69.89	87.35	47.91	90.19	56.06
IMP [2]	89.95	68.96	92.30	74.13	88.50	54.14	90.81	59.84
Proto-HATT [9]	91.54	70.26	93.43	75.24	90.63	57.26	92.86	61.51
Proto-AWATT [16]	93.35	75.37	95.28	80.16	92.06	65.65	93.42	69.70
LPN (o, o)	94.15	76.19	95.85	80.37	94.03	65.72	94.98	69.22
LPN (w, o)	96.41	82.26	97.43	85.81	95.26	71.25	96.23	75.49
LPN (w, w)	96.45	82.22	97.15	84.90	95.36	71.42	96.55	76.51

5.3 Implementation Details

Evaluation Metric. We follow [16] to adopt two widely used metrics Area Under Curve (AUC) and macro-F1 score to evaluate the performance.

Parameter Settings. For all experiments, we use the pre-trained language model Bert [7] to encode each word (token). Inspired by

[19], we freeze the first 6 layers of Bert and fine-tune the final 6 layers. For the model parameters, we set $d = 768$, $d' = 256$, $R = 4$ and $k = 100$ consistently. For the loss function, we set $\tau = 0.1$, $\lambda = 0.1$ and $\gamma = 0.01$ consistently and use AdamW [23] optimizer with the initial learning rate $1e-5$. For these parameters, we use the grid searching strategy and validation set to determine them. Take parameter k as an example. Figure 3 shows the performance of LPN

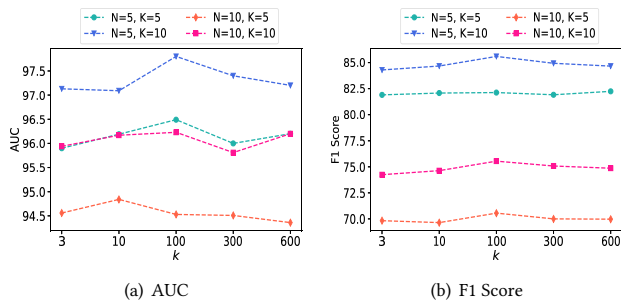


Figure 3: The performance of LPN (w,w) with different k values on the validation set of FewAsp.

(w,w) with different k values on the validation set of FewAsp. It can be seen that when k increases from 3 to 600, the performance improves at first and then drops, so we set $k = 100$ in experiments. Table 3 summarizes the main hyperparameters of our model.

5.4 Result Analysis

We perform experiments with 5/10-way and 5/10-shot settings on FewAsp (single), FewAsp (multi) and FewAsp three datasets. All reported results are from 5 different runs, and in each run the results are averaged over 600 test episodes. Table 4, 5 and 6 show the experimental results for FewAsp (single), FewAsp (multi) and FewAsp respectively. The baseline results are taken from [16] and the best results are highlighted in bold. From the results, we could make the following observations.

(1) LPN performs much better than other baselines. Specifically, in terms of AUC, LPN improves upon the most competitive baseline Proto-AWATT by 1.53%-2.13%, 2.66%-4.71% and 1.87%-3.30% on FewAsp (single), FewAsp (multi) and FewAsp respectively. In terms of macro-F1 score, LPN improves upon Proto-AWATT by 5.86%-9.12%, 5.11%-8.39% and 4.74%-6.85% on FewAsp (single), FewAsp (multi) and FewAsp respectively. The reason is that LPN leverages label description as auxiliary knowledge to learn more discriminative prototypes, integrates with contrastive learning to obtain better embeddings and uses a more effective multi-label inference module to accurately compute the aspect count.

(2) For all the methods, the results on FewAsp (multi) are a little worse than those on FewAsp (single) and FewAsp. The reason is that FewAsp (multi) consists of a large amount of sentences with multiple aspects, which increases the complexity of the dataset greatly. However, the proposed LPN can still achieve the best performance compared with other baselines, which further demonstrates the superiority of LPN in dealing with more complex multi-label tasks.

5.5 Ablation Study

Label-enhanced Prototypes. To verify the effectiveness of label-enhanced prototypes, we make the ablation study. The results are shown in Table 4, 5 and 6. LPN (o,o) means the LPN model without label enhancement and contrastive learning, and LPN (w,o) means the LPN model with label enhancement and without contrastive learning. It is easy to find that LPN (w,o) always performs much better than LPN (o,o) in all cases, which validates the effectiveness

of the label-enhanced prototypes. The reason is that label text descriptions contain lots of aspect-relevant semantic information, which is highly conducive to obtain more discriminative prototypes. **Contrastive Learning.** We also make the ablation study for the module of contrastive learning. The results are shown in Table 4, 5 and 6. LPN (w,o) means the LPN model with label enhancement and without contrastive learning, and LPN (w,w) means the LPN model with label enhancement and contrastive learning. We can observe that in most cases LPN (w,w) performs better than LPN (w,o). This is because that contrastive learning module can push samples in the same class close and samples in different classes further apart, thus obtaining better sentence embeddings.

5.6 Visualization

To better observe how the embeddings change with label-enhanced prototypes and contrastive learning, we sample 3000 episodes from test set of FewAsp (multi) in 5-way-5-shot setting, and then use t-SNE [34] to visualize the prototype embeddings obtained from LPN (o,o), LPN (w,o) and LPN (w,w). Note that we originally intend to visualize the sentence embeddings, but each sentence may contain multiple aspects which is difficult to distinguish by color. As each prototype is associated with unique aspect and is generated by the corresponding intra-class sentences, it can represent the sentence embedding to some extent. Figure 4 gives the visualization result of prototype embeddings obtained from LPN (o,o), LPN (w,o) and LPN (w,w). Prototypes (data points) with the same color contains the same aspect. It is easy to find that the distribution generated by LPN (o,o) has a lot of overlaps. The label enhancement in LPN (w,o) can help to separate the embeddings to some extent. The contrastive learning can further guarantee that the embeddings from same class are pulled together and the embeddings from different classes are pushed apart.

6 CONCLUSION

In this paper, we propose a label-enhanced prototypical network (LPN) to deal with multi-label few-shot aspect category detection. To learn more discriminative prototypes, LPN adopts the label text description as auxiliary knowledge to retain aspect-relevant information while eliminating the negative effect triggered by irrelevant aspects. To obtain better sentence embeddings to facilitate the aspect category detection task, LPN introduces contrastive learning to reduce intra-class discrepancy and enlarge the inter-class difference among sentence embeddings. Extensive experiments on three real-world datasets show that LPN outperforms the state-of-the-art methods by a large margin. In future work, we plan to investigate the theoretical underpinnings of our approach and extend our model to other multi-label few-shot scenarios like intent detection in dialogue systems.

ACKNOWLEDGMENTS

The authors are grateful to the anonymous reviewers for their valuable comments and suggestions. This work was supported by National Natural Science Foundation of China (No. 62106035, 61876028), and Fundamental Research Funds for the Central Universities (No. DUT20RC(3)040, DUT20RC(3)066).

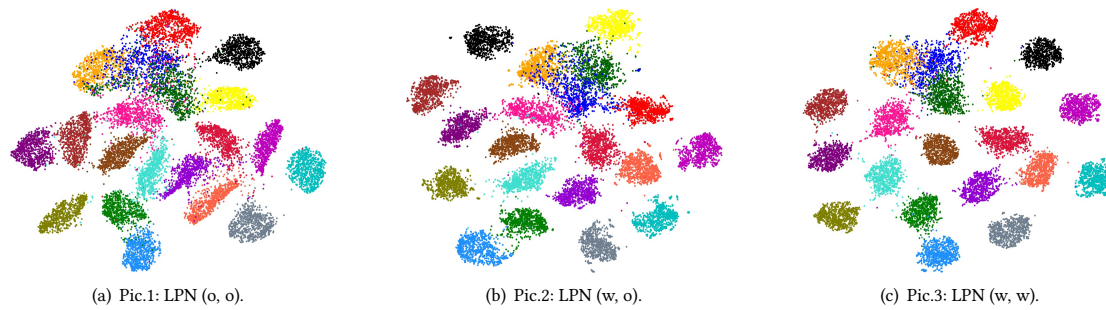


Figure 4: Visualization of prototype embeddings obtained from LPN (o, o), LPN (w, o) and LPN (w, w) respectively.

REFERENCES

- [1] Amit Alfassy, Leonid Karlinsky, Amit Aides, Joseph Shtok, Sivan Harary, Rogério Schmidt Feris, Raja Giryes, and Alexander M. Bronstein. 2019. LaSO: Label-Set Operations Networks for Multi-Label Few-Shot Learning. In *CVPR*. 6548–6557.
- [2] Kelsey R. Allen, Evan Shelhamer, Hanul Shin, and Joshua B. Tenenbaum. 2019. Infinite Mixture Prototypes for Few-shot Learning. In *ICML*. 232–241.
- [3] Yujia Bao, Menghua Wu, Shiyu Chang, and Regina Barzilay. 2020. Few-shot Text Classification with Distributional Signatures. In *ICLR*.
- [4] Konstantin Bauman, Bing Liu, and Alexander Tuzhilin. 2017. Aspect Based Recommendations: Recommending Items with the Most Valuable Aspects Based on User Reviews. In *KDD*. 717–725.
- [5] Hao Chen, Yaohui Wang, Benoit Lagadec, Antitza Dantcheva, and François Brémond. 2021. Joint Generative and Contrastive Learning for Unsupervised Person Re-Identification. In *CVPR*. 2004–2013.
- [6] Kai-Hsiang Cheng, Szu-Yu Chou, and Yi-Hsuan Yang. 2019. Multi-label Few-shot Learning for Sound Event Recognition. In *IEEE International Workshop on Multimedia Signal Processing (MMSp Workshop)*. 1–5.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT*. 4171–4186.
- [8] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In *ICML*, Vol. 70. 1126–1135.
- [9] Tianyu Gao, Xu Han, Zhiyuan Liu, and Maosong Sun. 2019. Hybrid Attention-Based Prototypical Networks for Noisy Few-Shot Relation Classification. In *AAAI*. 6407–6414.
- [10] Tianyu Gao, Xu Han, Hao Zhu, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2019. FewRel 2.0: Towards More Challenging Few-Shot Relation Classification. In *EMNLP*. 6249–6254.
- [11] Ruiying Geng, Binhua Li, Yongbin Li, Yuxiao Ye, Ping Jian, and Jian Sun. 2019. Few-Shot Text Classification with Induction Network. *CoRR* abs/1902.10482 (2019).
- [12] Erfan Ghadery, Sajad Movahedi, Hesham Faily, and Azadeh Shakery. 2019. MNCC: A Multilingual Ngram-Based Convolutional Network for Aspect Category Detection in Online Reviews. In *AAAI*. 6441–6448.
- [13] Beliz Gunel, Jingfei Du, Alexis Conneau, and Veselin Stoyanov. 2021. Supervised Contrastive Learning for Pre-trained Language Model Fine-tuning. In *ICLR*.
- [14] Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018. FewRel: A Large-Scale Supervised Few-shot Relation Classification Dataset with State-of-the-Art Evaluation. In *EMNLP*. 4803–4809.
- [15] Yutai Hou, Yongkui Lai, Yushan Wu, Wanxiang Che, and Ting Liu. 2021. Few-shot Learning for Multi-label Intent Detection. In *AAAI*. 13036–13044.
- [16] Mengting Hu, Shiwan Zhao, Honglei Guo, Chao Xue, Hang Gao, Tiegang Gao, Renhong Cheng, and Zhong Su. 2021. Multi-Label Few-Shot Learning for Aspect Category Detection. In *ACL*. 6330–6340.
- [17] Mengting Hu, Shiwan Zhao, Li Zhang, Keke Cai, Zhong Su, Renhong Cheng, and Xiaowei Shen. 2019. CAN: Constrained Attention Networks for Multi-Aspect Sentiment Analysis. In *EMNLP*. 4600–4609.
- [18] Manoj Kumar, Varun Kumar, Hadrien Glaude, Cyprien de Lichy, Aman Alok, and Rahul Gupta. 2021. Protoda: Efficient Transfer Learning for Few-Shot Intent Classification. In *IEEE Spoken Language Technology Workshop (SLT Workshop)*. 966–972.
- [19] Jaejun Lee, Raphael Tang, and Jimmy Lin. 2019. What Would Elsa Do? Freezing Layers During Transformer Fine-Tuning. *CoRR* abs/1911.03090 (2019).
- [20] Yuncong Li, Cunxiang Yin, Sheng-hua Zhong, and Xu Pan. 2020. Multi-Instance Multi-Label Learning Networks for Aspect-Category Sentiment Analysis. In *EMNLP*. 3550–3560.
- [21] Zhouhan Lin, Minwei Feng, Cícero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017. A Structured Self-Attentive Sentence Embedding. In *ICLR*.
- [22] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. 2019. Few-Shot Unsupervised Image-to-Image Translation. In *ICCV*. 10550–10559.
- [23] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *ICLR*.
- [24] Sajad Movahedi, Erfan Ghadery, Hesham Faily, and Azadeh Shakery. 2019. Aspect Category Detection via Topic-Attention Network. *CoRR* abs/1901.01183 (2019).
- [25] Tsensuren Munkhdalai and Hong Yu. 2017. Meta Networks. In *ICML*. 2554–2563.
- [26] Sachin Ravi and Hugo Larochelle. 2017. Optimization as a Model for Few-Shot Learning. In *ICLR*.
- [27] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy P. Lillicrap. 2016. One-shot Learning with Memory-Augmented Neural Networks. *CoRR* abs/1605.06065 (2016).
- [28] Victor Garcia Satorras and Joan Bruna Estrach. 2018. Few-Shot Learning with Graph Neural Networks. In *ICLR*.
- [29] Kim Schouten, Onne van der Weijde, Flavius Frasincar, and Rommert Dekker. 2018. Supervised and Unsupervised Aspect Category Detection for Sentiment Analysis with Co-occurrence Data. *IEEE Transactions on Cybernetics* 48, 4 (2018), 1263–1275.
- [30] Jake Snell, Kevin Swersky, and Richard S. Zemel. 2017. Prototypical Networks for Few-shot Learning. In *NeurIPS*. 4077–4087.
- [31] Qi Su, Kun Xiang, Houfeng Wang, Bin Sun, and Shiwen Yu. 2006. Using Pointwise Mutual Information to Identify Implicit Features in Customer Reviews. In *ICCPOL*. 22–30.
- [32] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip H. S. Torr, and Timothy M. Hospedales. 2018. Learning to Compare: Relation Network for Few-Shot Learning. In *CVPR*. 1199–1208.
- [33] Tun Thura Thet, Jin-Cheon Na, and Christopher S. G. Khoo. 2010. Aspect-based sentiment analysis of movie reviews on discussion boards. *Journal of Information Science* 36 (2010), 823–848.
- [34] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of Machine Learning Research* 9, 11 (2008).
- [35] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. 2016. Matching Networks for One Shot Learning. In *NeurIPS*. 3630–3638.
- [36] Peng Wang, Kai Han, Xiu-Shen Wei, Lei Zhang, and Lei Wang. 2021. Contrastive Learning Based Hybrid Networks for Long-Tailed Image Classification. In *CVPR*. 943–952.
- [37] Yikai Wang, Chengming Xu, Chen Liu, Li Zhang, and Yanwei Fu. 2020. Instance Credibility Inference for Few-Shot Learning. In *CVPR*. 12833–12842.
- [38] Guangfeng Yan, Lu Fan, Qimai Li, Han Liu, Xiaotong Zhang, Xiao-Ming Wu, and Albert Y. S. Lam. 2020. Unknown Intent Detection Using Gaussian Mixture Model with an Application to Zero-shot Intent Classification. In *ACL*. 1050–1060.
- [39] Jaesik Yoon, Taesup Kim, Ousmane Dia, Sungwoong Kim, Yoshua Bengio, and Sungjin Ahn. 2018. Bayesian Model-Agnostic Meta-Learning. In *NeurIPS*. 7343–7353.
- [40] Dian Yu, Luheng He, Yuan Zhang, Xinya Du, Panupong Pasupat, and Qi Li. 2021. Few-shot Intent Classification and Slot Filling with Retrieved Examples. In *NAACL-HLT*. 734–749.
- [41] Zhou Yu, Jun Yu, Chenchao Xiang, Jianping Fan, and Dacheng Tao. 2018. Beyond Bilinear: Generalized Multimodal Factorized High-Order Pooling for Visual Question Answering. *IEEE Transactions on Neural Networks and Learning Systems* 29, 12 (2018), 5947–5959.
- [42] Xinjie Zhou, Xiaojun Wan, and Jianguo Xiao. 2015. Representation Learning for Aspect Category Detection in Online Reviews. In *AAAI*. 417–424.