

Unsupervised Fusion Feature Matching for Data Bias in Uncertainty Active Learning

Wei Huang, Shuzhou Sun, Xiao Lin, Ping Li, *Member, IEEE*, Lei Zhu, Jihong Wang, C. L. Philip Chen, *Fellow, IEEE*, and Bin Sheng, *Member, IEEE*

Abstract—Active learning (AL) aims to sample the most valuable data for model improvement from the unlabeled pool. Traditional works, especially uncertainty-based methods, are prone to suffer from a data bias issue, which means that selected data cannot cover the entire unlabeled pool well. Although there have been lots of literature works focusing on this issue recently, they mainly benefit from the huge additional training costs and the artificially designed complex loss. The latter causes these methods to be redesigned when facing new models or tasks, which is very time-consuming and laborious. This paper proposes a feature matching-based uncertainty that resamples selected uncertainty data by feature matching, thus removing similar data to alleviate the data bias issue. To ensure that our proposed method does not introduce a lot of additional costs, we specially design an Unsupervised Fusion Feature Matching (UFFM), which does not require any training in our novel AL framework. Besides, we also redesign several classic uncertainty methods to be applied to more complex visual tasks. We conduct rigorous experiments on lots of standard benchmark datasets to validate our work. The experimental results show that our UFFM is better than the similar unsupervised feature matching technologies, and our proposed uncertainty calculation method outperforms random sampling, classic uncertainty approaches, and recent state-of-the-art uncertainty approaches.

Index Terms—Active learning, feature fusion, feature matching, neural network, uncertainty, data bias, deep learning.

Manuscript received April 13, 2021; revised June 29, 2022; accepted September 18, 2022. This work was supported in part by the National Natural Science Foundation of China under Grant 61872241 and Grant 61572316, in part by the National Key Research and Development Program of China under Grant 2019YFB1703600, and in part by The Hong Kong Polytechnic University under Grant P0030419, Grant P0042740, and Grant P0035358. (The first three authors contributed equally to this work.)

Wei Huang is with the School of Optical-Electrical and Computer Engineering, University of Shanghai for Science and Technology, Shanghai 200093, China (e-mail: 191380039@st.usst.edu.cn).

Shuzhou Sun is with the College of Information, Mechanical and Electrical Engineering, Shanghai Normal University, Shanghai 200234, China (e-mail: 1000479143@smail.shnu.edu.cn).

Xiao Lin is with the College of Information, Mechanical and Electrical Engineering, Shanghai Normal University, Shanghai 200234, China, and also with the Shanghai Engineering Research Center of Intelligent Education and Bigdata, Shanghai 200240, China (e-mail: lin6008@shnu.edu.cn).

Ping Li is with the Department of Computing and the School of Design, The Hong Kong Polytechnic University, Hong Kong (e-mail: p.li@polyu.edu.hk).

Lei Zhu is with the ROAS Thrust, The Hong Kong University of Science and Technology (Guangzhou), Guangzhou 511400, China, and also with the Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology, Hong Kong (e-mail: leizhu@ust.hk).

Jihong Wang is with the Shanghai University of Sport, Shanghai 200438, China (e-mail: cylwysy@163.com).

C. L. Philip Chen is with the School of Computer Science and Engineering, South China University of Technology, Guangzhou 510006, China, and also with the Pazhou Lab, Guangzhou 510335, China (e-mail: philip.chen@ieee.org).

Bin Sheng is with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: sheng-bin@sjtu.edu.cn).

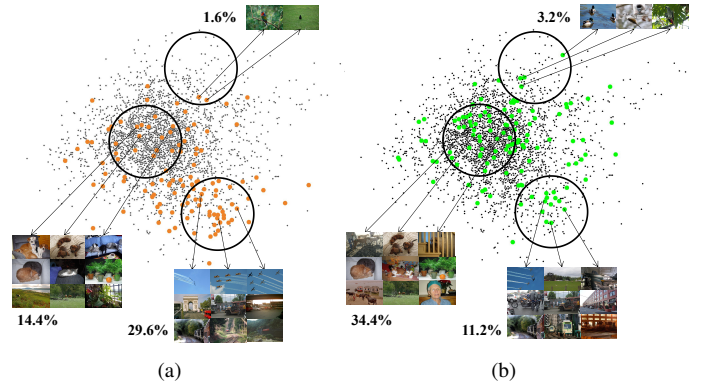


Fig. 1: t-SNE embedding of images on PASCAL VOC 2007 training set under the task of object detection. The black points are the distribution of the original data, and the orange and cyan points are the sampled data by the uncertainty-based active learning method. (a) is obtained from the baseline approach (Entropy Sampling [1]), and (b) is the sampling results of our method.

I. INTRODUCTION

IN this era of data flooding, labeling all of them for supervised learning is very time-consuming and laborious and thus is not realistic [2]–[5]. Albeit exhilarating prosperity in the line of semi-supervised [6] and unsupervised learning [7], however, supervised learning is still better than those two technologies in most scenarios [8]. Therefore, it is a key way to improve the model performance to take some of the most valuable data from the unlabeled pool for supervised learning, which is what Active Learning (AL) does [9], [10]. Existing active learning methods are mainly divided into three groups: 1) Uncertainty-based approaches [3], [9], [11]–[14]. This kind of algorithm calculated the uncertainty of unlabeled data based on the current learned model. 2) Diversity-based approaches [13], [15]. Diversity approaches preferentially selected the batch of data with the most dispersed feature distance. 3) Expected model change approaches [10], [16]. They took the processed unlabeled data (e.g., adding noise) as the inputs to observe the changes of the model outputs. In addition to these methods, some recent works also considered the relationship between these sampling strategies. Multi-Criteria Active Deep Learning [17] selected informative samples by considering multiple criteria simultaneously (i.e., density, similarity, etc.), and it has achieved excellent performance on the classification tasks. In this paper, our focus is to develop an uncertainty

active learning framework. Compared to the other two types of approaches, the uncertainty-based method has a lower cost when facing large-scale unlabeled data.

The uncertainty active learning approach is to select the data with the worst confidence, which is similar to the leak filling in the human learning process. However, because the proportion and difficulty of categories in the training data are different, the learned model will inevitably tend to go overboard on partial categories. For example, if the current learned model has a poor learning effect on the *bird* category, it will give a high uncertainty to all data containing birds in the unlabeled pool, which is the data bias problem [18]. Fig. 1 illustrates an example of the data bias problem. We use t-distributed Stochastic Neighbor Embedding (t-SNE) [19] to show the distribution of data to be sampled (black dots), and we use orange and cyan dots to represent the sampled data. The baseline approach (i.e., Entropy Sampling (ES) [1]) sampled a large amount of data in the sparse area (the bottom black circle area) of the unlabeled pool but selected a few samples in the dense area (the middle black circle area). However, Our proposed active learning framework resamples the uncertainty data obtained by original uncertainty approaches, and our sampling results are obviously more able to cover the entire unlabeled data pool.

Confidence estimation is a common manner to calculate the uncertainty, and traditional uncertainty sampling methods, including Least Confidence [8], [9], Margin Sampling [20], Entropy Sampling [1], [2], etc. also belong to this. However, when the model to be learned is deep neural networks (DNNs), the probability distributions obtained by the traditional uncertainty approaches are too confident, which will lead to the serious data bias problem and thus be even worse than random sampling [18]. To alleviate the data bias problem, recent literature has sought improvement from multiple perspectives. Wasserstein Adversarial Active Learning (WAAL) [21] adopted a Wasserstein distance to refactor the uncertainty calculation method to alleviate the data bias. Ensembles-based active learning (ENS) [3] used an ensemble network to calculate data uncertainty. Loss Prediction Module (LPM) [8] took the unlabeled data as a part of model training to predict target losses of unlabeled inputs. However, these methods will introduce additional calculation and training costs. Also, they are all highly task-related, which means that they need to be redesigned when facing other tasks.

Unsupervised feature matching does not require any training resources or unbearable computational costs. And it can be used to calculate the similarity between unlabeled data [22]–[24]. This fact motivates us to utilize feature-matching to compute the similarity in selected uncertainty data for alleviating the data bias problem in uncertainty AL approaches. To achieve this goal, we first propose Unsupervised Fusion Feature Matching (UFFM), which can calculate the data similarity from the perspective of multi-layer network features. Then, we design a novel uncertainty calculation method, which resamples uncertainty data obtained by other basic uncertainty active learning methods and then removes the redundant ones, thereby significantly alleviating the data bias problem. Our method can be combined with any current uncertainty

approaches to improve their performance.

Last but not least, considering that Least Confidence, Margin Sampling, Entropy Sampling, etc., are only applicable to classification task, we have also redesigned those methods to make them suitable for the object detection task. In summary, the contributions of our work are three-fold:

- We propose an efficient active learning framework, which is to resample the selected uncertainty data based on feature matching to alleviate the problem of data bias. Compared with other methods that focus on this problem, our proposed method has a lower cost and can be combined with all existing uncertainty methods to improve their performance.
- We design Unsupervised Fusion Feature Matching (UFFM), which fuses multiple layers of features to generate descriptors for feature matching. Our approach can perceive the details of the feature more comprehensively, while other similar methods can only perceive very limited information.
- We improve several uncertainty methods originally designed for classification tasks such that these uncertainty estimation manners can be adapted for handling complex images in the object detection task. The reason behind is that our uncertainty computation method considers all objects in the image to calculate their uncertainty, thereby providing a more reliable uncertainty estimation than original uncertainty methods.

II. RELATED WORK

Here, we first review the most related works about deep learning-based feature matching, which can be roughly classified as supervised-based and unsupervised-based approaches. Then, we present existing active learning methods.

A. Deep learning-based feature matching

Early descriptors are often hand-designed [25], [26]. Recently, many researchers focused on developing deep learning-based methods for learning features due to their impressive performance in diverse vision tasks. Here we mainly review deep learning-based approaches, including supervised approaches and unsupervised approaches.

Supervised approaches. To solve the overfitting issue caused by a lack of training data, HashGAN [27] synthesized nearly real images to augment the training set and thus could obtain high-quality descriptors for image matching. Deep Spherical Quantization (DSQ) [28] used a CNN to generate supervised and compact descriptors for image matching. Meanwhile, DSQ forced the network to leverage a L_2 normalization to alleviate the negative effect of norm variance. Deep Product Quantization (DPQ) [29] introduced a dictionary-based representation to ensure a more accurate image matching and classification under maintaining an affordable computational complexity and memory. Shen et al. [30] added two additional fully-connected layers onto the top of the backbone network to obtain binary descriptors, and it showed that a simple network can also extract effective semantic information

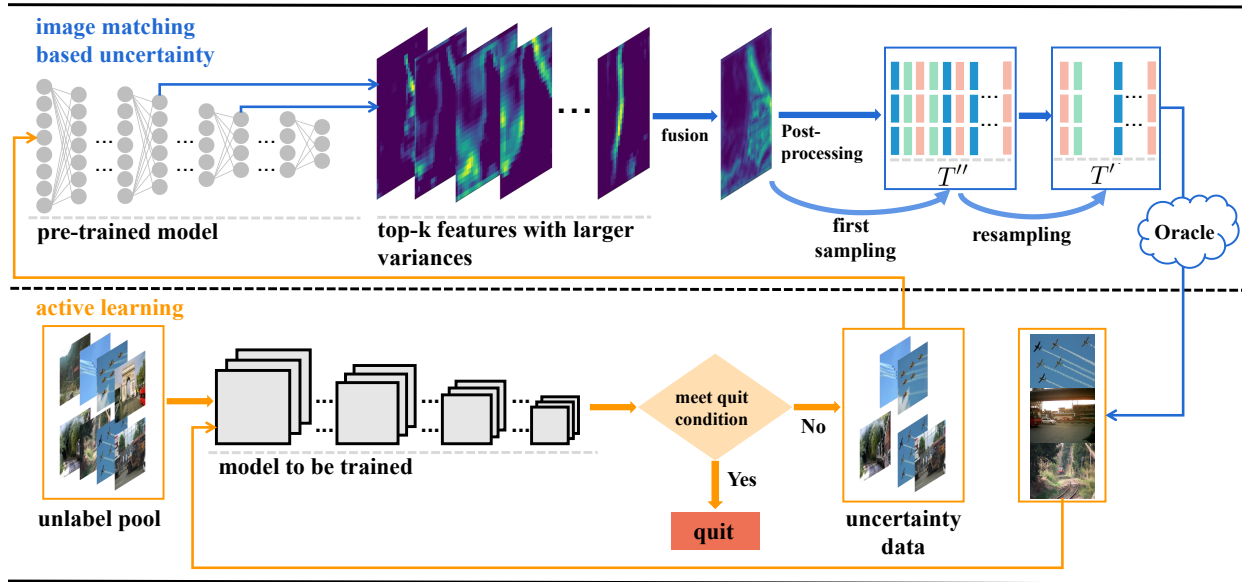


Fig. 2: The active learning framework consisting of Unsupervised Fusion Feature Matching (UFFM) and feature matching based uncertainty. Our framework first selects a certain amount of uncertainty data (more than the expected of active learning) from the unlabeled pool through a basic uncertainty approach. Then, we use UFFM to resample the above-selected data to remove the similar data. The quit condition is the label budget is exhausted or the expected performance of the model is reached.

through a proper design. Although supervised-based approaches have achieved remarkable performance, collecting large-scale labeled datasets to train the feature matching network is a challenging task, especially for specific fields [25]. Meanwhile, supervised-based approaches require prohibitive training costs but still have an overfitting risk.

Unsupervised approaches. Unsupervised feature matching has drawn widespread attention in recent years due to its low costs on training data. Similarity-Adaptive Deep Hashing (SADH) [31] alternatively proceeded with three training modules (i.e., deep hash model training, similarity graph updating, and binary code optimization), which helped to update the similarity graph matrix more effectively than traditional methods. DistillHash [32] obtained the descriptor through the relationship between the initial signals learned from local structures and the semantic similarity labels assigned by the optimal Bayesian classifier. Deep variational binaries (DVB) [33] introduced the conditional auto-encoding variational Bayesian networks to exploit the feature space structure of the training data using the latent variables better to unveil the intrinsic structure of the whole sample space. However, these unsupervised methods often require a tedious redesign when facing different models, and it is difficult to guarantee that the intrinsic structure of the whole sample space can be obtained when facing different data sets. Instead, Part-based Weighting Aggregation (PWA) [34] proposed a pure unsupervised feature matching method, which directly aggregated the features of the pre-trained model. However, from their experiments, we find that this type of method cannot fully perceive the features when the difference between the pre-training dataset and the images to be matched was large. For addressing this problem, in this work, we propose to fuse the middle layer and the lower

layer features to get a more complete descriptor for images.

B. Active learning approaches

Uncertainty-based active learning approaches. As one of the most commonly-used methods in active learning, uncertainty-based methods are prone to data bias problems when facing large-scale data or DNNs. Apart from the above classical approaches, some recent advanced methods have achieved different performance improvements from multiple perspectives. Modeling Active Learning (SMAL) [11] combined uncertainty, diversity, and density via sparse modeling to alleviate the data bias problem. However, the sparse representation is challenging to guarantee stability when facing large-scale data. Batch Mode Active Learning (BMAL) [12] started with a feature descriptor extraction coupled with a divergence matrix to alleviate the problem of redundancy between unlabeled points. However, the traditional feature extraction method used in BMAL can hardly contribute to DNNs. Loss Prediction Module (LPM) [8] was jointly trained with the target model to predict the target loss of unlabeled inputs, but it increases the costs of network training. Localization-Aware Active Learning (L-Aware) [4] proposed a localization tightness and localization stability to calculate the uncertainty. However, this work requires the network to provide intermediate prediction results (e.g., predictions by Region Proposal Network (RPN) in Faster R-CNN [35]), which means that this method cannot be used in a model without intermediate prediction (e.g., one-stage object detector [36]–[38]). Ensemble-based method [3] used five committee networks to calculate uncertainty, which tends to be impractical in the existence of large-scale unlabeled data and deep neural networks. In this paper, we use a feature matching

algorithm to resample uncertain data obtained by uncertainty AL methods. Our work is based on a pre-trained model and thus does not introduce any training costs. Meanwhile, our method can be combined with any existing uncertainty AL methods and thus has an excellent generalization ability.

Diversity-based active learning approaches. Diversity-based approaches preferred to select the batch of data with the most dispersed feature distance. Patra and Bruzzone [13] used a kernel k -means clustering algorithm to minimize the redundancy and keep the diversity among these samples after selecting a batch of uncertain samples. Yang et al. [39] regarded active learning as a discrete optimization problem, and they imposed a diversity constraint on the objective function to make the selected data as diverse as possible. The Core-set approach [40] improved the competitiveness of selected data by constructing a core subset. Although these methods have shown to be effective for simple and low-scale features, our empirical analysis suggests that they do not scale to learn more complex and large-scale features. We will compare this type of method and prove our point in the experiment.

Expected model change active learning approaches. Expected model change approaches take processed unlabeled data (e.g., adding noise) as inputs to observe the changes of the model outputs. Settles et al. [41] estimated the value of unlabeled data by measuring the changes of model parameters, but it ignored the underlying data distribution. To address this issue, Freytag et al. [42] directly calculated the expected change of model predictions and marginalized the unknown label. Furthermore, Käding et al. [43] proposed a new generalization of the Expected Model Output Change principle and thus this expected model change active learning approaches can be used in DNNs. However, compared with uncertainty-based methods, this kind of method has a higher cost, especially when faced with DNNs and large-scale unlabeled data.

III. OUR METHOD

This section presents the implementation details of our proposed framework. First, we will introduce our Unsupervised Fusion Feature Matching (UFFM), which can remove similar data in the sampling results obtained by uncertainty active learning approaches. Then, we will propose a novel uncertainty calculation technology coupled with UFFM to calculate the uncertainty of unlabeled data. The former can obtain the data with higher uncertainty, while the latter can further eliminate the data bias in the above-selected data. We finally introduce the implementation details of the proposed active learning framework based on special tasks, including the classification task and the object detection task. Fig. 2 shows the schematic illustration of our proposed deep active learning method.

A. Unsupervised Fusion Feature Matching (UFFM)

Selection of layers. For DNN models, the lower layers often detect the surface features of objects (e.g., edges, textures, shapes, etc.), while the higher layers reflect more abstract information (e.g., classification, etc.), and this has also been concluded in many works [4], [34]. Unlike existing methods that only use the high layer features of the pre-trained model

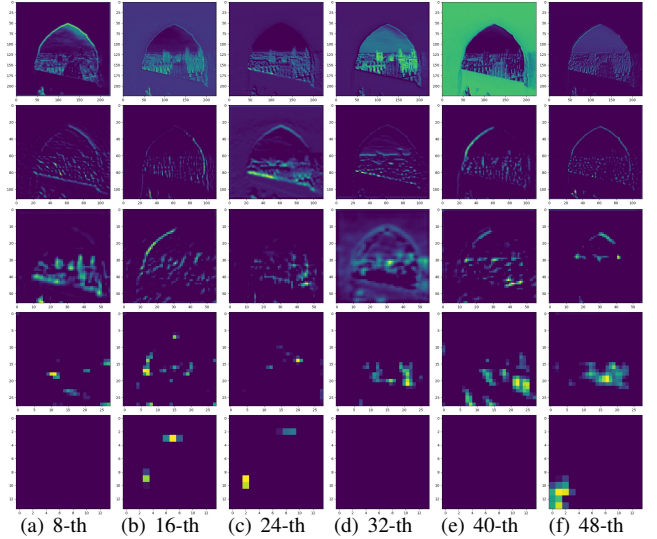


Fig. 3: The feature visualization of different channels in different layers. From top to bottom are block1_conv2, block2_conv2, block3_conv3, block4_conv3, and block5_conv3. The used model is VGG16 pre-trained on ImageNet, and the test image is selected from Oxford5k.

for unsupervised feature matching, we choose features at the lower and middle layers in this paper. Specifically, for the l -layers model $\mathcal{L}_{i=1}^l$, we select the i' -th layer $\mathcal{L}_{i'}$ and the i'' -th layer $\mathcal{L}_{i''}$ for unsupervised feature matching, where $1 \leq i' \leq i'' \leq l$, i' and i'' are determined according to the used model.

Such selection of layers is mainly based on two reasons. The first reason is that the high layer of the pre-trained model tends to provide very limited useful features, because the pre-training dataset and the images to be matched may be very different, and the pre-trained model cannot recognize the abstract features of the images to be matched. To prove our point, we show the inference results of images in Oxford5k through the model pre-trained by ImageNet [44], where the similarity between Oxford5k and ImageNet is very limited. Fig. 3 shows that the high layer of the pre-trained model can hardly detect the abstract features of the image to be matched. Another reason is that the neural network has information loss during the downsampling, and thus the selection of both the lower and middle layers at the same time helps to obtain the complete features of the image to be matched. We will further discuss the benefits of fusing different layers in more detail in the ablation study.

Selection of channels. For the above-selected layers, we only choose partial channels for feature matching. It has three main considerations: 1) Feature maps of DNN models usually have many channels, and the direct use of all these channels for feature matching undoubtedly requires a very large computational cost. 2) As high-dimensional features, more channels are more like to contain noise, which is obviously not expected for the feature matching task. 3) Much previous literature [34] argues that channels with larger variances are more discriminative. To verify this view, we visualize the channels with different variances, as shown in Fig. 4. From

these visualization results, we can find that channels with larger variances can often detect fixed information (see from two to four columns), while channels with small variances can hardly reflect any features (see the last three columns). Therefore, this paper only uses channels with larger variances for feature extraction. For the selected layer \mathcal{L}_l , we suppose that its shape is $h_l \times w_l \times c_l$, where c_l is the number of channels, and h_l and w_l are the height and width of channels. We first calculate the variances of all channels $\mathcal{V}_l = v_1, v_2, \dots, v_{c_l}$ in \mathcal{L}_l , where the i -th channel variance v_i can be calculated as:

$$v_i = \frac{1}{m \times n} \sum_{m=1}^{h_l} \sum_{n=1}^{w_l} (x_{(m,n)}^i - \bar{x}^i)^2, \quad (1)$$

where $x_{(m,n)}^i$ is the value of position (m,n) at the i -th channel. \bar{x}^i is the average variance of the i -th channel, and it can be calculated by:

$$\bar{x}^i = \frac{1}{m \times n} \sum_{m=1}^{h_l} \sum_{n=1}^{w_l} (x_{(m,n)}^i). \quad (2)$$

Then we sort the variances of all channels in descending order and select the top k channels for unsupervised feature matching. k is determined by the used model, and we will discuss it in more detail in the ablation study.

Feature fusion. Since a single layer cannot extract features completely, our UFFM uses the lower and middle layer channels to describe the features. In recent years, many works have used the features of convolutional layers for feature matching. Compared with fully-connected layers, the convolutional layers is more interpretable. However, instead of using channels of a single convolutional layer, we also choose that of the lower layer to ensure the completion of the image description. Pooling layers are pooled from the convolutional layer, so it can keep the features of the convolutional layer. Meanwhile, the relationship between the lower pooling layer and the higher convolutional layer of the current deep neural network is often multiple, so the fusion of these two types of layers will be reasonable and simple.

For the selected lower pooling layer $\mathcal{L}_{l'}$ and the middle convolutional layer $\mathcal{L}_{l''}$, we suppose that their shapes are $h_{l'} \times w_{l'} \times c_{l'}$ and $h_{l''} \times w_{l''} \times c_{l''}$ respectively, where $h_{l'}/h_{l''} = w_{l'}/w_{l''} = N$, $N \in (1, 2, 3, \dots)$. For these two selected layers, we first pool $\mathcal{L}_{l'}$ to the same size as $\mathcal{L}_{l''}$. Then we select the corresponding meaningful channels from the above layers and aggregate them. SPoC [23] argued that the aggregation of the convolution layer channels could be directly based on a sum pooling aggregation, without the need to use a fusion method like Fisher Vector and triangular embedding, since the convolutional layer features of deep neural networks were sufficiently discriminative. This paper also does not use the traditional aggregation method. We adopt the approach of a bitwise addition to aggregate $\mathcal{L}_{l'}$ and $\mathcal{L}_{l''}$ into the shape of $h_{l''} \times w_{l''} \times (c_{l'} + c_{l''})$. We show the details of feature fusion in Fig. 5. We will further prove the effectiveness of this design in ablation study, especially the selection of fused features.

Post-processing of fused features. For the image I , we can get its fused features $f_I \in \mathbb{R}^{h_{l''} \times w_{l''} \times (c_{l'} + c_{l''})}$ through above steps. We first perform a l_2 -normalization on f_I to

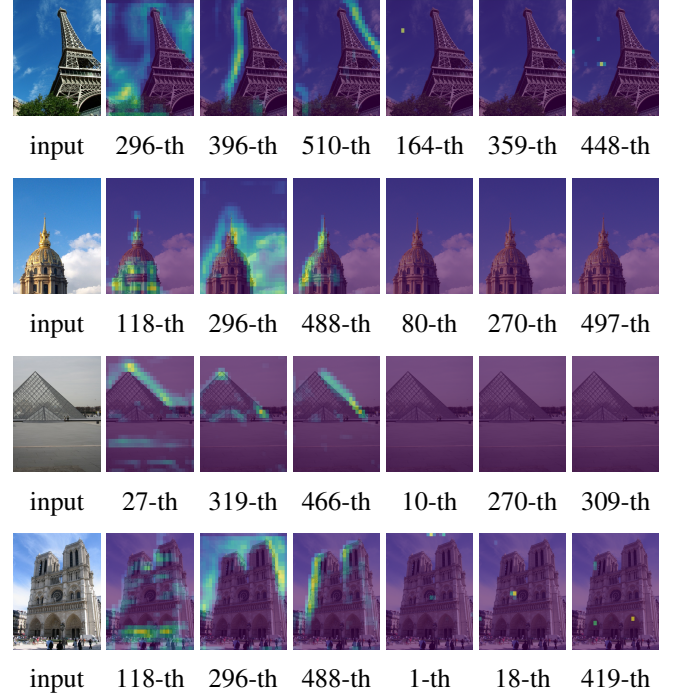


Fig. 4: The channels with different variances. Two to four columns are the three channels with the largest variance, and the last three columns are the ones with the smallest variance. The pre-trained model used here is VGG16 trained by ImageNet, and input images are selected from Oxford5k.

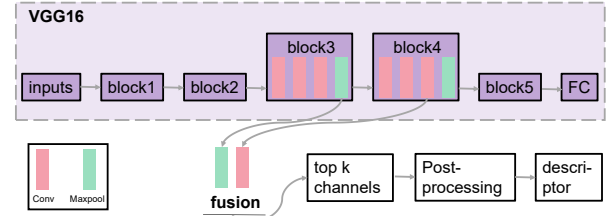


Fig. 5: The details of feature fusion. The network we used in this paper is VGG16, and the fused features are block3_pool and block4_conv3.

obtain a new feature map, $f_I' = \frac{f_I}{\|f_I\|_2}$. Then, we use PCA (Principal Component Analysis) to reduce the dimension of the normalized features, and denote the result f_I , $f_I \in \mathbb{R}^{h_{l''} \times w_{l''} \times (c_{l'} + c_{l''})'}$, where $h_{l''}'$, $w_{l''}'$ and $(c_{l'} + c_{l''})'$ can be adjusted according to the expected complexity.

B. Statement for data bias in active learning

Let $Q_{(x,y)}$ and $P_{(x,y)}$ denote the distribution of unlabeled data pool and the selected data obtained by an AL method, and suppose their densities are $q(x,y) = q(y | x)q(x)$ and $p(x,y) = p(y | x)p(x)$, respectively. We use $\mathcal{H}(h \sim H)$ to represent the optimal sampling for the original distribution H under the condition of a given sampling rate, where h obey the distribution H . Based on this definition, $\mathcal{H}((x,y) \sim P_{(x,y)})$ can be calculated as:

$$\mathcal{H}((x,y) \sim Q_{(x,y)}) = - \iint q(y | x)q(x) \ln(q(y | x)q(x)) dx dy. \quad (3)$$

Algorithm 1 Feature Matching Based Uncertainty**Input:** $D_{t=1}^T, T', T'', (T' < T'')$ **Output:** $S_{t=1}^{T''}$

- 1: Compute the uncertainty of $D_{t=1}^T$, and select T'' data,
 $\mathcal{D}_{t=1}^{T''} = \mathcal{U}(U_{(D_{t=1}^T)}, T'')$;
- 2: $q = 0$
- 3: **for** $m = 1$ to T'' **do**
- 4: Add d_m to the set of $S_{t=1}^{T''}$;
- 5: Compute the similarity,
 $M_{d_m} = UFFM(d_m, d_{(m+1) \sim T''})$;
 Mark the data in $d_{(m+1) \sim T''}$ as similar or dissimilar
 according to M_{d_m} ;
- 6: **for** $n = m + 1$ to T'' **do**
- 7: **if** d_m and d_n are dissimilar **then**
- 8: Add d_n to the set of $S_{t=1}^{T''}$;
- 9: $q = q + 1$;
- 10: **if** $q \geq T'$ **then**
- 11: break;
- 12: **end for**
- 13: **end for**

$$\mathcal{H}((x, y) \sim P_{(x, y)}) = - \int \int q(y | x) q(x) \ln(p(y | x) p(x)) d_x d_y. \quad (4)$$

We then use KL divergence $D_{KL}(Q_{(x, y)} \parallel P_{(x, y)})$ to describe the extent to which $P_{(x, y)}$ covers $Q_{(x, y)}$:

$$D_{KL}(Q_{(x, y)} \parallel P_{(x, y)}) = \mathcal{H}((x, y) \sim P_{(x, y)}) - \mathcal{H}((x, y) \sim Q_{(x, y)}) = \int \int q(y | x) q(x) \ln \frac{q(y | x) q(x)}{p(y | x) p(x)} d_x d_y. \quad (5)$$

Therefore, we can obtain the optimal active learning query function \mathcal{Q}_{AL} by minimizing $D_{KL}(Q_{(x, y)} \parallel P_{(x, y)})$:

$$\mathcal{Q}_{AL} = \arg \min_{P_{(x, y)}} D_{KL}(Q_{(x, y)} \parallel P_{(x, y)}). \quad (6)$$

However, from an example shown in Fig. 1, we can see that $P_{(x, y)}$ is biased towards partial categories in practice. Assuming that the optimal sampling of training data under given conditions \mathcal{Q}_{AL} . Obviously, in $P_{(x, y)}$, some high uncertainty data $\mathcal{Q}_{AL} \setminus \mathcal{H}((x, y) \sim P_{(x, y)})$ are not queried by \mathcal{Q}_{AL} , but some low uncertainty data $\mathcal{H}((x, y) \sim Q_{(x, y)}) \setminus \mathcal{Q}_{AL}$ are selected instead.

C. Feature matching based uncertainty

The existing uncertainty based active learning approaches are prone to suffer from the data bias problem since the learned models often have a preference for partial data. To alleviate this problem, we propose a novel uncertainty method, which improves the original uncertainty approach via feature matching to get the uncertainty of unlabeled data. Specifically, we first use the original uncertainty approach to get the uncertainty of the unlabeled data. Then, we use the UFFM proposed in this paper to resample the selected data by the original approach to alleviate the data bias problem.

Let $D_{t=1}^T$ denote the unlabeled dataset, and T is the amount of data. Then, the original active learning method selects fixed number of data to support the model training:

$$\mathcal{D}_{t=1}^{T'} = \mathcal{U}(U_{(D_{t=1}^T)}, T'), \quad (7)$$

where U is the original uncertainty approach, which can calculate the uncertainty of unlabeled data based on current learned network. \mathcal{U} means to select T' data from $D_{t=1}^T$ according to the uncertainty. $\mathcal{D}_{t=1}^{T'} = \{d_1, d_2, \dots, d_{T'}\}$ is selected data, and T' is the expected number of data to be selected at the current stage. Unlike original approaches, our proposed method first select T'' data from the unlabeled pool:

$$\mathcal{D}_{t=1}^{T''} = \mathcal{U}(U_{(D_{t=1}^T)}, T''), \quad (8)$$

where $T'' > T'$. $\mathcal{D}_{t=1}^{T''} = \{d_1, d_2, \dots, d_{T''}\}$ denotes the selected T'' unlabeled data. All data in $\mathcal{D}_{t=1}^{T''}$ are sorted according to their uncertainty scores. Taking the j -th data d_j in $\mathcal{D}_{t=1}^{T''}$ as an example, we calculate the similarity through UFFM:

$$M_{d_j} = UFFM(d_j, d_{(j+1) \sim T''}), \quad (9)$$

where $d_{(j+1) \sim T''} = [d_{j+1}, d_{j+2}, \dots, d_{T''}]$. $UFFM(d_j, d_{(j+1) \sim T''})$ can calculate the similarity between $d_j \in \mathbb{R}^{1 \times \mathcal{D}_s}$ and $d_{(j+1) \sim T''} \in \mathbb{R}^{(T''-j) \times \mathcal{D}_s}$, where \mathcal{D}_s is the dimension of the descriptor. Specifically, we obtain the distance through a matrix multiplication $d_j \times [d_{(j+1) \sim T''}]^\top$, and the shape of result is $\mathbb{R}^{1 \times (T''-j)}$. We binarize the similarity by marking 10% data in $d_{(j+1) \sim T''}$ with the smallest distance as similar and the others as dissimilar. Lastly, we add the data that is not similar to d_j to the AL results $S_{t=1}^{T''}$. For all data in $\mathcal{D}_{t=1}^{T''}$, we perform the above steps in turn until the amount of data reaches the expected number T' . We use pseudo code to show our proposed sampling strategy, see Algorithm 1 for more details.

In addition, for an actual application (e.g., image classification or object detection) of active learning, T' is often a fixed value. However, T'' is a hyperparameter in our proposed framework. Here we define the Sampling Rate (SR), $SR = \frac{T''}{T'}$. Obviously, $SR=1$ means that our framework degenerates to the original uncertainty methods. However, if SR is large, the uncertainty of the data selected by our method may be lower. We set $SR = 1.2$ in this paper (i.e., $T'' = 1.2 \times T'$), and the influence of SR has been further discussed in ablation study.

D. Further design

We will verify our proposed method on image classification and object detection. Because these two tasks cover classification and regression, the generalization of our proposed method can be fully illustrated.

Further design for image classification. Least Confidence (LC), Margin Sampling (MS), and Entropy Sampling (ES) are classic uncertainty methods for the image classification task. Least Confidence calculates the uncertainty of unlabeled data through the maximum predicted probability, while Margin Sampling and Entropy Sampling consider the first two and all probabilities, respectively. For the image I in the dataset with c classes, Least Confidence uncertainty I_{LC} , Margin Sampling uncertainty I_{MS} , Entropy Sampling uncertainty I_{ES} can be calculated as follows:

$$I_{LC} = (1 - \hat{C}), \text{ s.t. } \hat{C} = \arg \max_{i \in [1, \dots, c]} (p_i), \quad (10)$$

$$I_{MS} = (\hat{C} - \hat{C}'), s.t. \hat{C} = \arg \max_{i \in [1, \dots, c]} (p_i), \hat{C}' = \arg \max_{i \in [1, \dots, c] \setminus \hat{C}} (p_i), \quad (11)$$

$$I_{ES} = \sum_{i=1}^c p_i \log(p_i), \quad (12)$$

where p_i represents the confidence of the i -th class. Further, we can use Algorithm 1 to obtain feature matching based uncertainty.

Further design for object detection. The above three classic uncertainty methods do not consider the situation that an image may contain multiple objects to be recognized in object detection. Hence, we redesign these three methods to make them more suitable for object detection. Our redesigned methods include Redesigned Least Confidence (RLC), Redesigned Margin Sampling (RMS), and Redesigned Entropy Sampling (RES), and their uncertainties I_{RLC} , I_{RMS} and I_{RES} can be calculated as follows:

$$I_{RLC} = \sum_{j=1}^{n_I} (1 - \hat{C}_j), s.t. \hat{C}_j = \arg \max_{i \in [1, \dots, c]} (p_i^j), \quad (13)$$

$$I_{RMS} = \sum_{j=1}^{n_I} (\hat{C}_j - \hat{C}_j'), \quad (14)$$

$$s.t. \hat{C}_j = \arg \max_{i \in [1, \dots, c]} (p_i^j), \hat{C}_j' = \arg \max_{i \in [1, \dots, c] \setminus \hat{C}_j} (p_i^j),$$

$$I_{RES} = \sum_{j=1}^{n_I} \sum_{i=1}^c p_i^j \log(p_i^j), \quad (15)$$

where n_I is the number of objects in image I . p_i^j represents the confidence that the j -th prediction box in the image is i -th class. Similarly, we use these redesigned methods coupled with Algorithm 1 to obtain feature matching based uncertainty. Apart from the above three classic uncertainty methods, we also compare the state-of-the-art uncertainty approaches, see experiments for more.

IV. EXPERIMENTAL RESULTS

In this paper, we propose a method that can alleviate the data bias problem in active learning. To verify the effectiveness of our work, we have conducted many experiments. First, we introduce the datasets for feature matching, classification, and object detection. Next, we verify our proposed UFFM on the task of feature matching. Then, we prove that our proposed active learning framework can achieve competitive performance on the task of classification and detection. Finally, we discuss the design details of the proposed framework and its advantages in mitigating data bias through the ablation study.

A. Datasets

1) *Datasets for feature matching:* **Oxford5k** [52]. The Oxford buildings dataset consists of 5,026 images collected from Flickr by searching particular Oxford landmarks. This dataset includes 11 different landmarks, and each landmark



Fig. 6: Matching results by UFFM. The query images and results in the first two lines are from Oxford, and the other lines are from Paris. The pre-trained model here is VGG16, and the selected fusion features are block3_pool and block4_conv3.

TABLE I: Performance comparison between pure unsupervised feature matching approaches (P) and fine-tuning based image matching approaches (F) under different feature descriptors (d). The pre-trained model here is VGG16, and the selected fusion features are block3_pool and block4_conv3.

method	d	Oxford5k		Paris6k	
		P	F	P	F
NetVLAD [24]	128	55.5	63.5	64.3	73.5
MAC [45]	128	55.7	76.8	70.6	78.8
ours	128	56.2	-	70.9	-
SPoC [23]	256	53.1	-	-	-
R-MAC [45]	256	56.1	78.2	72.9	83.5
RVD-W [46]	256	60	-	-	-
Razavian et al. [22]	256	-	67	-	53.3
ours	256	62.5	-	73.3	-
CroW [47]	512	70.8	-	79.7	-
InterActive [48]	512	65.6	-	79.2	-
PWA [34]	512	72	87.8	82.3	94.9
CNNBoW [45]	512	-	79.7	-	83.8
ours	512	73.2	-	84.5	-

contains 5 query images with ground truth. For landmarks in each image, it contains one of the following possible labels: i) Good. It means that the landmark is very clear. ii) OK. It represents that at least 25% of the landmark is clear. iii) Bad. This means the landmark is not present. iv) Junk. It represents that less than 25% of the landmark is visible. **Paris6k** [53]. The Paris dataset consists of 6,412 images collected from Flickr by searching Paris landmarks. The label format of this dataset is the same as the Oxford5k. Both Oxford5k and Paris6k are standard datasets for evaluating feature matching methods.

2) *Dataset for classification:* **CIFAR-10** [54]. The CIFAR-10 dataset consists of 60k color images with 10 categories, (6k images per category). The dataset has 50k training images and 10k test images. **Fashion-MNIST** [55]. The Fashion-MNIST is a fashion product dataset, including 60k training images and 10k test images. The dataset has 10 categories, which are more

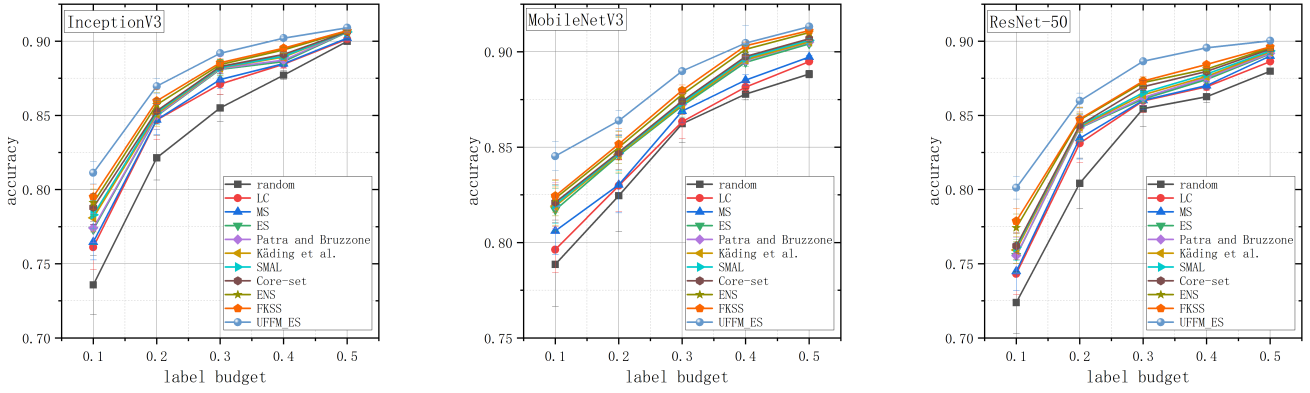


Fig. 7: Comparisons under Metric 1 (i.e., performance under fixed labeling budget) on CIFAR 10. The pre-trained model we used here is VGG16 trained by ImageNet. We repeat each experiment five times and report standard deviation by error bar. The evaluation networks used here are InceptionV3 [49], ResNet-50 [50], and MobileNetV3 [51]. The compared baselines included classic uncertainty methods (LC [9], MS [20], ES [1]), recent uncertainty-based methods (ENS [3], Patra and Bruzzone [13], SMAL [11], and FKSS [18]), recent diversity-based methods (Core-set [40]), recent expected model change methods (Kading et al. [43]), and random sampling. The basic uncertainty approach combined with UFFM is Entropy Sampling (ES).

difficult than the original MNIST dataset. Both CIFAR-10 and Fashion-MNIST are classic classification datasets.

3) *Dataset for object detection: PASCAL VOC 2007* [56]. VOC 2007 contains 20 object categories, and it includes 2.5k training images, 2.5k validation images, and 5k test images. **PASCAL VOC 2012**. VOC 2012 is an augmented version of VOC 2007, which contains about 5k training images and 5k validation images. VOC 2007 and VOC 2012 are both standard datasets commonly used for vision tasks, including classification, detection, segmentation, etc.

B. Evaluation of our proposed Unsupervised Fusion Feature Matching

Our Unsupervised Fusion Feature Matching (UFFM) is a feature matching technology, which can cooperate with original uncertainty approaches to calculate the uncertainty of the unlabeled data to alleviate the data bias problem. We follow the evaluation protocol of Oxford and Paris to crop the image with the provided bounding box. The matching results by our UFFM are shown in Fig. 6.

Meanwhile, we report the quantitative matching results under the metric of mAP in Table I. Here we mainly compare two lines of methods using pure unsupervised feature matching [23], [24], [34], [45], [47] and starting with unsupervised coupled with fine-tuning [22], [34], [45]. From the quantitative results in Table I, we have the following observations: 1) Our UFFM outperforms existing pure unsupervised feature matching in all descriptor dimensions. We argue that this mainly benefits from reasonable channel selection and fusion, and we will prove this point in the ablation study. 2) Although our method is slightly inferior to fine-tuning-based approaches, our training and matching costs are far less. Meanwhile, considering that our method is mainly designed for the uncertainty calculation of active learning, we believe that this small gap does not significantly affect active learning. 3) Generally, fine-tuning-based methods outperform pure unsupervised methods.

TABLE II: Results under Metric 2 (i.e., labeling budget under expected performance) on Fashion-MNIST. Similarly, we repeat each experiment five times and report the average budget.

method	InceptionV3 [49]			ResNet-50 [50]			MobileNetV3 [51]		
	0.85	0.9	0.95	0.85	0.9	0.95	0.85	0.9	0.95
random	3.4	4.4	5.2	4.2	5.2	5.8	3.6	4.6	5.2
LC [9]	3.0	4.0	4.8	3.8	5.0	5.6	3.2	4.0	4.8
MS [20]	2.8	3.8	4.8	3.6	4.8	5.4	3.0	3.8	4.8
ES [1]	2.8	3.6	4.6	3.6	4.6	5.4	3.0	3.6	4.8
Kading et al. [43]	2.8	3.6	4.4	3.6	4.4	5.4	2.8	3.6	4.6
SMAL [11]	2.8	3.4	4.4	3.6	4.4	5.4	2.8	3.6	4.4
Core-set [40]	2.6	3.4	4.4	3.4	4.4	5.4	2.8	3.6	4.4
ENS [3]	2.6	3.4	4.2	3.4	4.4	5.2	2.8	3.6	4.2
FKSS [18]	2.6	3.4	4.2	3.2	4.2	5.2	2.6	3.4	4.2
ours	2.0	3.0	4.0	3.0	4.0	5.0	2.0	3.0	4.0

It shows that the pre-trained model cannot fully sense the features of the image to be matched, especially the abstract features. This is consistent with the conclusion we observed in Fig. 3. Therefore, this further illustrates the rationality of our abandonment of the high layer that can only provide limited features. 4) We argue that our method is better than other methods in helping to alleviate the data bias problem, because our method performs more complete feature matching through feature fusion, and these are exactly the features that the model wants to learn. That is, our method will examine the feature similarity among unlabeled data more comprehensively when calculating the uncertainty. Instead, other approaches only match the very limited features provided by the high layer.

C. Evaluation of our active learning method on the image classification task

Experimental setting. We explore three classic classification networks including InceptionV3 [49], ResNet-50 [50], and MobileNetV3 [51] as evaluation models of the active learning.

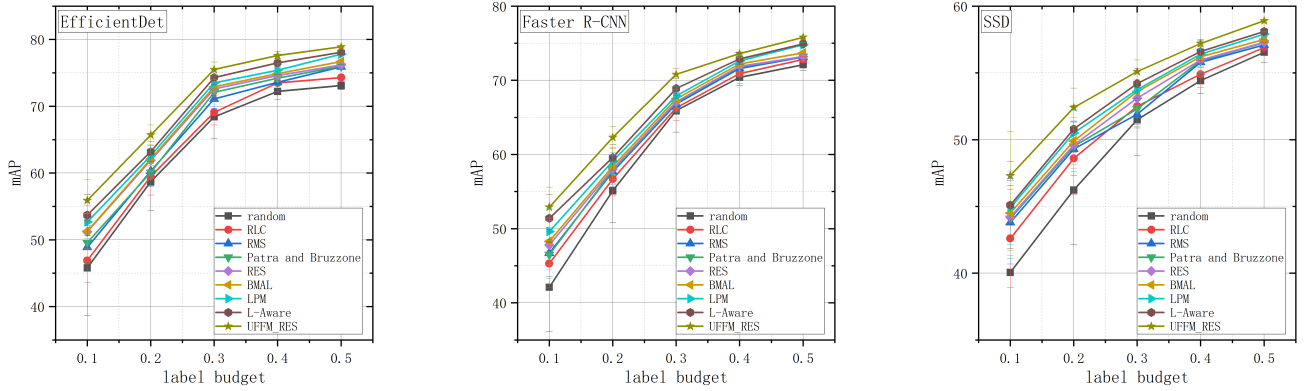


Fig. 8: Performance on PASCAL VOC 2007 under fixed labeling budget. The pre-trained model we used here is VGG16 trained by ImageNet. We repeat each experiment five times and report standard deviation by error bar. The evaluation detectors used here are EfficientDet [36], Faster R-CNN [35], and SSD [37]. The compared baselines included classic uncertainty methods (RLC [9], RMS [20], and RES [1]), recent state-of-the-art methods (Patra and Bruzzone [13], L-Aware [4], BMAL [12], and LPM [8]), and random sampling. The basic uncertainty approach combined with UFFM is Redesigned Entropy Sampling (RES).

TABLE III: Labeling budget under expected performance. The datasets we used here are PASCAL VOC 2007 and VOC 2012. Similarly, we repeat each experiment five times and report the average budget.

method	EfficientDet [36]			Faster R-CNN [35]			SSD [37]		
	0.7	0.75	0.8	0.7	0.75	0.8	0.55	0.60	0.65
random	2.8	4.0	5.2	3.0	4.0	5.2	3.2	4.2	5.2
RLC [9]	2.4	3.6	4.6	2.8	3.6	4.6	2.8	3.6	4.8
RMS [20]	2.4	3.4	4.4	2.6	3.6	4.6	2.6	3.6	4.8
RES [1]	2.4	3.4	4.4	2.6	3.6	4.4	2.4	3.6	4.6
BMAL [12]	2.4	3.4	4.4	2.4	3.6	4.4	2.4	3.4	4.6
LPM [8]	2.4	3.2	4.2	2.4	3.4	4.4	2.4	3.4	4.4
L-Aware [4]	2.2	3.2	4.2	2.4	3.2	4.2	2.2	3.2	4.4
ours	2.0	3.0	4.0	2.0	3.0	4.0	2.0	3.0	4.0

For all three classification networks, their hyper-parameters are: optimizer is SGD, weight_decay = 0.00004, decay factor of learning rate = 0.94, learning_rate = 0.01, momentum = 0.9, and batch size = 32. The basic settings of active learning are: the pre-trained model is VGG16, and the descriptor length is 256. We evaluate different active learning methods using the following two metrics:

Metric 1: performance under fixed labeling budget. A larger score under Metric 1 indicates a better active learning method. If the selected data requires a massive amount of labeling costs, this active learning algorithm loses its meaning. Therefore, we set the labeling budget for all active learning frameworks to be less than 50% of the original unlabeled data for fair comparisons.

Metric 2: labeling budget under expected model performance. The less the labeling costs are, the better the sampling approaches are. Hence, a smaller score under Metric 2 indicates a better active learning framework. The expected performance should be adjusted according to the model used, and please refer to Table II for details.

Compared methods. We compare our proposed method

against the following baseline methods: 1) Random sampling: sampling the data uniformly at random from the unlabeled set. 2) Classical uncertainty approaches: Least Confidence (LC) [9], Margin Sampling (MS) [20], and Entropy Sampling (ES) [1]. 3) Recent SOTA methods: Patra and Bruzzone [13] used a kernel k -means clustering algorithm to minimize the redundancy and kept the diversity among these samples after selecting a batch of uncertain samples. To highlight the advantages of our method in alleviating data bias, the baseline method [13] used the same uncertainty calculation and feature extraction mechanism as our method. The only difference is that Patra and Bruzzone [13] obtained the final sampling results through the clustering algorithm, while our proposed method uses a feature matching algorithm. Note that unlike the original paper [13], we use k -means++ as the clustering algorithm and run it multiple times to get better performance. Ensembles-based active learning (ENS) [3] used an ensemble network to calculate data uncertainty. Note that the training data in the first stage of this original paper is generated on the basis of labels, and the purpose is that the initial sets are balanced over all classes. But the real raw unlabeled data will not have any label information, so our framework is purely random sampling in the first epoch. For the sake of fairness, the data of the first epoch of ENS in this paper is also a purely random sampling. Fisher Kernel Self Supervision (FKSS) [18] proposed a low-complexity feature density matching method and utilized it to calculate the uncertainty of unlabeled data. Similarly, for fairness, we use the complete standard dataset when comparing this method, rather than using only partial data of the dataset in the original paper to create artificial data. The Core-set approach [40] is a diversity-based active learning technology, and we follow the training tricks and hyperparameters in the original paper. Kading et al. [43] is an expected model change AL technology. Following the paper [43], we also use a stochastic gradient approximation with just a single sample to estimate model parameter updates, and the models we use are all DNNs to ensure that the baseline

can play its advantages, thereby ensuring a fair comparison. Sparse Modeling Active Learning (SMAL) [11] combined uncertainty, diversity, and density via sparse modeling to alleviate the data bias problem. SMAL [11] divided the dataset into a seed set (labeled set), an unlabeled set, a validation set, and a testing set. For a fair comparison, we use the randomly selected data in the first epoch as the seed set and then follow the original paper to set other details of this method.

Results and analysis. We conducted experiments on CIFAR 10 and Fashion-MNIST with the above two metrics, including Metric 1 and Metric 2. Our proposed framework and other compared baseline methods follow the same training process. We first randomly sample 5% of the dataset (about 3k images) as the training data for the first epoch. Then for each subsequent epoch, we use the active learning framework to sample 5% of the dataset and use it to continue training the model. Finally, for Metric 1, we repeat active learning sampling and training until the sampled data reaches the fixed labeling budget. For Metric 2, we repeat until the trained model reaches the expected performance. The basic uncertainty approach combined with our proposed method UFFM here is Entropy Sampling (ES). We report the result on CIFAR 10 and Fashion-MNIST in Fig. 7 and Table II, respectively. And from these results, we have several observations:

1) Our method outperforms all baselines by a clear margin. Specifically, our framework has a higher performance under a fixed labeling budget (see Fig. 7), and our framework requires less labeled data under the expected performance (see Table II).

2) Considering that training data at the first stage is randomly sampled, we repeat each experiment five times and report the standard deviation (error bar in Fig. 7) for comparisons. From the results, we can clearly find that our method is more stable than compared methods. Note that the stability of active learning is related to the degree of data bias in each epoch. Hence, it indicates that our framework is superior to other methods in alleviating the data bias problem. We will prove this in detail in the ablation study.

3) As the fixed labeling budget or expected model performance grows, the improvement of our method over compared methods tends to decrease, as shown in Fig. 7. However, the application scenarios of active learning often have a small labeling cost. Hence, we believe that our method should have a superior performance in real cases.

4) Under part expected performance, there will be the epoch gap between our proposed framework and comparing baseline, our approach uses fewer epochs to achieve sample performance (e.g., our proposed AL framework with MobileNetV3 uses 4.0 epochs while the baseline of random sampling needs 5.2 epochs to meet the expected performance of 0.95). The epoch gap shows that our method can use less labeling budget, which means that more training resources can be saved.

D. Evaluation of our active learning method on the object detection task

Experimental setting. Here, we consider multiple detectors, including EfficientDet [36], Faster R-CNN [35], and

TABLE IV: Performance with 256-dimensional descriptors under different layer selections. The labeling budget used here is 0.5, and the experiment follows the settings in Fig. 7 and Fig. 8. The pre-trained model used here is VGG16 trained on ImageNet. Lower layer, middle layer, and high layer are block3_pool, block4_conv3, and block5_pool, respectively, and our method is the fusion of lower and middle layers.

method	feature matching		classification	detection			
	Oxford 5k	Paris 6k	ResNet-50	EfficientDet (D0)			
	mAP (%)	mAP (%)	accuracy (%)	mAP (%)	AP_S	AP_M	AP_L
lower layer only	60.9	71.5	88.6	77.8	40.9	82.1	89.1
middle layer only	61.6	72.1	89.1	78.1	43.3	82.7	89.5
high layer only	60.2	70.1	87.4	76.9	25.1	81.6	88.7
ours	62.5	73.3	90.0	78.9	44.9	83.6	90.3

SSD [37]. The settings of EfficientDet are: the backbone is EfficientDet-D0, the optimizer is SGD, initial learning rate=0.08, the warmup learning rate= 0.001, warmup steps =2500, and batch size = 128. The hyperparameters of Faster R-CNN are: the backbone is ResNet-101 [50], the optimizer is SGD, weight_decay = 0.00005, learning_rate = 0.0001, and batch_size = 32. The settings of SSD include: the backbone is MobileNetV2 [57], momentum is 0.94, and batchsize = 24. Moreover, the metrics used in the image classification task are also adopted in this section.

Compared methods. The following methods are employed for comparisons: 1) Random sampling. 2) Redesigned classical uncertainty approaches: RLC, RMS, and RES. 3) Recent SOTA methods: Patra and Bruzzone [13] is a diversity-based AL approach, and the comparison details are described in the previous section. Loss Prediction Module (LPM) [8] trained with the target active model, and it could be used to predict the target loss of unlabeled inputs. We use the same module and model connected to three layers of the target model for all the detectors. The internal structure and module training follow the settings in the original paper. Localization-Aware (L-Aware) [4] AL method used the localization tightness and the localization stability to calculate uncertainty. For Faster R-CNN, we use the region proposals provided by its RPN to calculate Localization tightness. While for EfficientDet and SSD, we directly calculate the localization stability since they do not have an intermediate proposal. BMAL [12] was a batch mode AL technique, which started with a feature descriptor extraction coupled with a divergence matrix to alleviate the problem of redundancy among unlabeled points. We follow most of the experimental details of BMAL [12]. For example, the Gabor filter is applied to the images for feature extraction, and PCA is used to reduce dimensionality. To reduce the computational costs, we follow the original paper to utilize a sub-sampling strategy. Also, we use BMAL [12] combined with the method we proposed in Section III to make it more suitable for object detection.

Results and analysis. We conduct solid evaluations on our active learning framework in terms of the object detection task and utilize two metrics (i.e., Metric 1 and Metric 2) proposed in Section 4.3 for comparisons. For Metric 1, we regard the training set of PASCAL VOC 2007 (about 2.5k images) as the

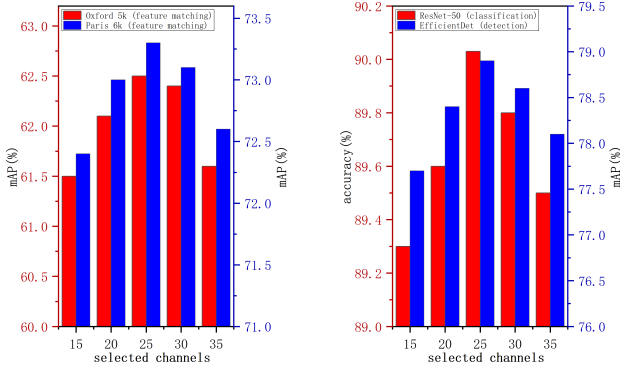


Fig. 9: Matching results with 256-dimensional descriptor under different selected top-k channels with maximum variances. The labeling budget used here is 0.5, and the experiment follows the settings in the above sections.

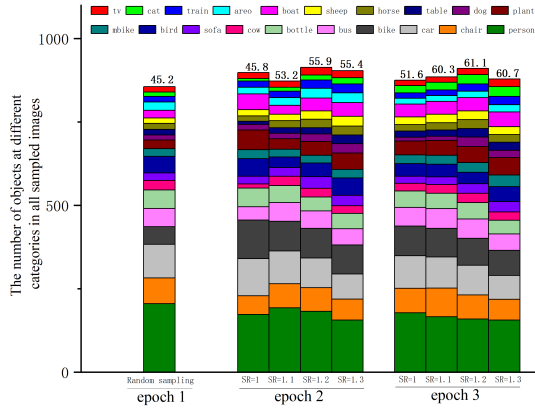


Fig. 10: The y-axis denotes the number of objects at each category in all sampled images under a specific epoch number. Here, in each bar, we also show the corresponding classification accuracy (e.g., 55.9% at the epoch 2 with $SR = 1.2$) of our method under a specific epoch number and a specific sampling rate (SR). The labeling budget used here is 0.5, and the experiment follows the settings in Fig. 7 and Fig. 8. The dataset and learned model are PASCAL VOC 2007 and EfficientDet, respectively. The pre-trained model for feature matching is VGG16, and block3_pool and block4_conv3 are selected features.

original unlabeled data $VOC_I^{2.5k}$. We let 5% of $VOC_I^{2.5k}$ as the training data for the first epoch, and in each subsequent epoch, we select 5% of the original unlabeled data as the new training data. The above step is repeated until the labeling budget is exhausted. While for Metric 2, we unite the train set of PASCAL VOC 2007 (about 2.5k images) and PASCAL VOC 2012 (about 5k images) as the original unlabeled data $VOC_I^{7.5k}$. We let 5% of $VOC_I^{7.5k}$ as the training set for the first epoch, and in each subsequent epoch, we select 5% of the original unlabeled data as the added training data. The above step is repeated until the trained model reaches the expected performance. The results of Metric 1 and Metric 2 are reported in Fig. 8 and Table III, respectively. The basic uncertainty approach combined with our proposed method UFFM here is

Redesigned Entropy Sampling (RES). From the above results, we can conclude that:

1) Our method outperforms all baselines on the object detection task. Meanwhile, object detectors we used in experiments cover two-stage (Faster R-CNN) and one-stage (EfficientDet and SSD), which shows that our framework is model-agnostic and thus can assist networks with any structure perform active learning.

2) From the standard deviation reported in Fig. 8, we find that our method is more stable than all baselines. We will demonstrate in the ablation study that this advantage is mainly due to the fact that our method can alleviate the data bias problem.

3) As the fixed labeling budget or the expected model performance increases, some uncertainty methods are even worse than random sampling. The main reason is that the problem of data bias has reached a serious degree. When the sampling epochs increase, the gap between our method and other methods is also decreased. However, we think this is normal. As active learning progresses, there is less and less valuable data, so the scope for model improvement will be limited.

4) The epoch gap of training in the image classification task can also be observed here.

E. Ablation study

Better layer combination. UFFM in this paper aims to fuse the features of different layers to perceive richer information. To further demonstrate the effectiveness of this fusion mechanism, we report the performance under different layer combinations in Table IV, and we find that our proposed method outperforms other approaches, especially for small objects. It demonstrates the advantages of our UFFM method in detailed feature perception. Meanwhile, we can also find that our method is the best, in which the middle layer is better than the lower layer, and the high layer is the worst. The main reason is that although the lower layers can perceive more features, it also introduces lots of noise. Since the pre-training dataset (ImageNet) is different from the target dataset (Oxford5k, Paris6k, CIFAR 10, PASCAL VOC, etc.), the perceptible features of the high layer are also limited. Note that the results in Table IV are not to deny the advantages of high layer features in the supervised task. It only shows that in the scene where the intersection of the images to be processed and the pre-training data is small, the high layer features are weaker than the middle and lower layer features.

Suitable channel selection. Many previous works have found that the perception ability of a specific channel is relatively fixed, and channels with larger variances perceive more information. In Fig. 9, we report the results under the top-k channels with the largest variances. We find that too small channel numbers cannot perceive complete information, while too larger channel numbers introduce additional noise. Based on the experimental results in Fig. 9, we empirically set the number of channels as 25.

More optimized sampling distribution. In Fig. 10, we report the sampling distribution of different object categories

TABLE V: The performance under redesigned uncertainty approaches and original uncertainty approaches. The dataset used here is PASCAL VOC 2007. OUA and RUA are original uncertainty approaches and redesigned uncertainty approaches, respectively.

	method	backbone	labeling budget			
			0.2	0.3	0.4	0.5
OUA	EfficientDet [36]	D0	62.5	73.8	76.2	77.1
	Faster R-CNN [35]	ResNet-101	59.8	68.9	72.7	74.6
	SSD [37]	MobileNetV2	51.1	54.2	56.9	58.2
RUA	EfficientDet [36]	D0	65.7	75.5	77.6	78.9
	Faster R-CNN [35]	ResNet-101	62.3	70.8	73.6	75.8
	SSD [37]	MobileNetV2	52.4	55.1	57.2	58.9

TABLE VI: The performance and training costs of different training techniques. Non-pretrain and pre-train stand for training from scratch and pre-training on ImageNet, respectively. The labeling budget used here is 0.5, and the experiment follows the settings in Fig. 7 and Fig. 8. Note that the pre-trained models here down from the Tensorflow GitHub repository, so the pre-training costs are not included in the training costs.

		classification (ResNet-50)		detection (EfficientDet (D0))	
		accuracy (%)	train cost (h)	mAP (%)	train cost (h)
Unsupervised (UFFM)		90.0	0.8	78.9	7.5
Supervised	Non-pretrained	90.5	1.3	79.3	14.0
	Pre-trained	91.3	1.1	79.8	12.5

and the model performance at different epochs under multiple Sampling Rates (SR). Sampling distribution can help us easily compare the changes of category quantity in different epochs, and the model performance can help us choose a better Sampling Rate. From Fig. 10, we observe that the original uncertainty method ($SR=1$) has an obvious data bias problem. For example, compared with the first epoch, the second epoch has a significant increase in the number of samples for bike, plant, boat, etc. According to the principle of the uncertainty active learning method, we think that it is because the current model has a poor learning effect on these categories. However, such a direct sampling method ignores the similarity of the internal data of the same classification, which causes the data bias problem. For the same uncertainty unlabeled data, the sampling results of our method obviously do not have the above problems, which means that the increase in the number of samples for bike, plant, boat, etc., does not fluctuate as much as the original method. This is because our proposed method can remove similar images in these image classifications. According to the results reported in Fig. 10, we set up $SR=1.2$ in this paper.

To further quantify the advantages of our work in mitigating the data bias problem, we define the coverage rate here, which can be used to measure the coverage degree of the sampled data to the original data. The coverage rate is S/O , where O is the original data and S is the residual data after removing the similar parts from the sampled data. We report the coverage rate in Fig. 11, which shows that our proposed active learning framework yields better sampling results than other methods.

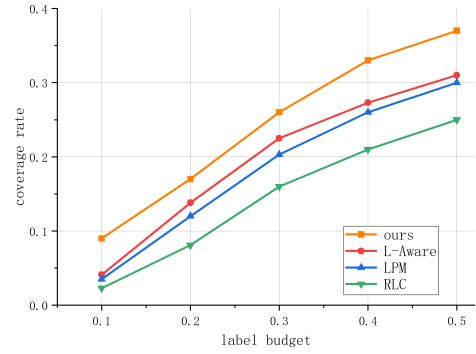


Fig. 11: Coverage rate under fixed labeled budget. The data and detector used here are EfficientDet and PASCAL VOC 2007, respectively, and the settings of the experiment are the same as that in Section 4.4. Feature matching uses UFFM proposed in this paper. The compared baselines include classic uncertainty methods RLC [9], recent state-of-the-art methods L-Aware [4] and LPM [8].

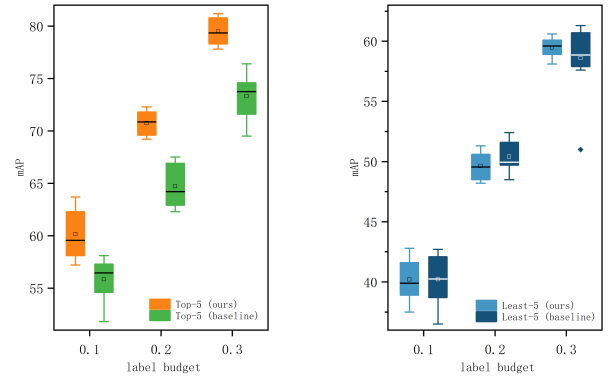


Fig. 12: Performance on PASCAL VOC 2007 under fixed labeling budget. The test set has 20 classes, Top-5 represents 5 classes with the highest mAP, and Least-5 represents 5 classes with the worst performance. Each experiment is repeated 10 times, and the baseline method is Entropy Sampling (ES).

More stable training process. Data bias causes an instability problem in the training process, which can also be drawn from Fig. 7 and Fig. 8. Here we will further study the underlying causes of this phenomenon. We think that this is mainly because the data bias incurs that partial categories do not progress steadily. To verify this, we report the learning process in Fig. 12, and we observe that for the baseline method, the progress of the Top-5 categories in the next epoch is very limited. This is mainly because the current learned model will focus on the Least-5 categories, which means that the next epoch will be mainly sampled Least-5 categories. Instead, our method will resample the sampling results of the Least-5 categories that are currently focused on. Because what we remove is the data that has a limited help for the model learning, we can guarantee the learning of the Top-5 categories without reducing the learning effect of the Least-5 categories, thereby ensuring a stable learning process.

More effective redesigned uncertainty approaches. We redesign several uncertainty-based approaches to make them

fit the object detection task. Those redesigned approaches consider the uncertainty of all objects in the image rather than just one object. This improvement is mainly because we have observed that the complexity of the images used for object detection is much greater than image classification, that is, the images used for detection often contain many objects, and the test confidence of these objects is different greatly because of the impact of category, size, etc. We report the performance of our proposed and original uncertainty approaches in Table V, and we find that our redesigned uncertainty approaches are more suitable for the task of object detection.

F. Feature Redundancy

Feature Redundancy in existing unsupervised feature matching. Unsupervised feature matching has attracted lots of attention in recent years. Using pre-trained models for feature extraction of images to be matched has almost become a de facto approach. ImageNet is a very large dataset, and its commonly-used subset ILSVRC2012 (ImageNet2012) still has more than one million training images. Currently, many vision models are pre-trained on ImageNet, and these pre-trained models used in unsupervised feature matching are also usually based on it. Oxford5k and Paris6k are standard datasets used for feature matching. Currently, there are many unsupervised feature matching works that use pre-trained models on ImageNet to evaluate the proposed methods on these two datasets.

However, we find in the experiment that the above-unsupervised feature matching paradigm has the Feature Redundancy problem. Feature Redundancy is defined here that non-matching targets in the images to be matched are perceived by the pre-trained model, which will have a negative impact on feature matching because the features of non-matching targets will interfere with feature matching. We argue that Feature Redundancy exists because of the category difference between the pre-training dataset and that of used for matching. For example, ILSVRC2012 has 1,000 categories, while Oxford5k and Paris6k have only a very limited number of categories. We illustrate an example of this problem in Fig. 13. People are the non-matching target in Paris6k, but it is a category of the pre-training dataset ILSVRC2012. Therefore, the pre-trained model can perceive the features of the human category, and this type of feature is the so-called “redundant features”. Redundant features negatively affect image matching because the dataset to be matched does not consider itself to include such features. That is, two similar images may be classified into different categories due to redundant features. Obviously, the more the category difference between the pre-training dataset and that of used for matching, the higher the possibility of feature redundancy.

The influence of feature redundancy. To further quantify the influence of feature redundancy, we compare our unsupervised method against supervised methods. In addition, to show the advantages of our proposed unsupervised method in reducing training costs, we compare two supervised methods, which are pre-trained on ImageNet and trained from scratch, respectively. We report the results in Table VI, and we have

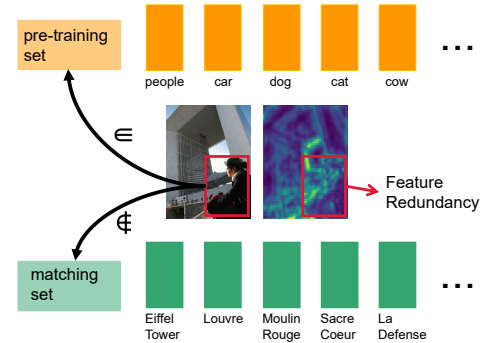


Fig. 13: An instance of Feature Redundancy. The test image is selected from Paris6k, the feature shown here is the fusion of the top-25 channels with the largest variance in block3_conv3 of VGG16.

several observations: 1) Feature redundancy has an impact on the model performance. The reason is that the supervised methods match the features of the target dataset more accurately, so it can better alleviate the data bias. 2) Our proposed unsupervised method saves considerable training costs at the expense of acceptable model performance (saving up to 38.5% and 46.4% of the training costs on the classification and detection). It is conceivable that in the face of large-scale data, our method will have greater advantages in cost saving. 3) Unlike supervised methods, our primary goal is to pursue a trade-off between the model performance and the training costs. Therefore, our method has a greater practical value, especially in the face of large-scale unlabeled data.

V. CONCLUSION

This paper presents a novel uncertainty calculation method to alleviate the problem of data bias in uncertainty-based active learning. By using Unsupervised Fusion Feature Matching (UFFM) to resample the selected uncertainty data, the feature matching method we design does not introduce too many additional costs. Our active learning framework based on feature matching outperforms random sampling, classic uncertainty approaches, and recent state-of-the-art uncertainty approaches in the task of image classification and object detection. Meanwhile, unlike those active learning methods that can only be used based on specific tasks or models, our framework is task-agnostic and model-agnostic and thus can be combined with almost any current uncertainty method to improve their performance. We have proved the effectiveness of our framework on image classification and object detection through experiments. In fact, our method can be applied to more complex vision tasks such as pedestrian re-identification and segmentation, and we take them as future work.

REFERENCES

- [1] A. J. Joshi, F. Porikli, and N. Papanikolopoulos, “Multi-class active learning for image classification,” in *Computer Vision and Pattern Recognition*, 2009, pp. 2372–2379.
- [2] H. H. Aghdam, A. Gonzalez-Garcia, A. Lpez, and J. Weijer, “Active learning for deep detection neural networks,” in *International Conference on Computer Vision*, 2019, pp. 3671–3679.

- [3] W. H. Beluch, T. Genewein, A. Nrnberger, and J. M. Khler, "The power of ensembles for active learning in image classification," in *Computer Vision and Pattern Recognition*, 2018, pp. 9368–9377.
- [4] C. C. Kao, T. Y. Lee, P. Sen, and M. Y. Liu, "Localization-aware active learning for object detection," in *Asian Conference on Computer Vision*, 2019, pp. 506–522.
- [5] L. Nie, M. Liu, and X. Song, *Multimodal Learning Toward Micro-Video Understanding*. Morgan & Claypool Publishers, 2019.
- [6] M. mjeja, M. Woczyk, J. Tabor, and B. C. Geiger, "SeGMA: Semi-supervised gaussian mixture autoencoder," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–12, 2020.
- [7] T. Ergen and S. S. Kozat, "Unsupervised anomaly detection with LSTM neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 8, pp. 3127–3141, 2020.
- [8] D. Yoo and I. S. Kweon, "Learning loss for active learning," in *Computer Vision and Pattern Recognition*, 2019, pp. 93–102.
- [9] D. D. Lewis and W. A. Gale, "A sequential algorithm for training text classifiers," in *Conference on Research and Development in Information Retrieval*, 1994, p. 312.
- [10] N. Roy and A. McCallum, "Toward optimal active learning through sampling estimation of error reduction," in *International Conference on Machine Learning*, 2001, p. 441448.
- [11] G. Wang, J.-N. Hwang, C. Rose, and F. Wallace, "Uncertainty-based active learning via sparse modeling for image classification," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 316–329, 2019.
- [12] S. Chakraborty, V. Balasubramanian, Q. Sun, S. Panchanathan, and J. Ye, "Active batch selection via convex relaxations with guaranteed solution bounds," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 10, pp. 1945–1958, 2015.
- [13] S. Patra and L. Bruzzone, "A cluster-assumption based batch mode active learning technique," *Pattern Recognition Letters*, vol. 33, no. 9, pp. 1042–1048, 2012.
- [14] Y. Lin, Y. Zhang, J. Chen, S. Zhang, and D. Z. Chen, "Suggestive annotation: A deep active learning framework for biomedical image segmentation," *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 399–407, 2017.
- [15] Z. Wang and J. Ye, "Querying discriminative and representative samples for batch mode active learning," in *International Conference on Knowledge Discovery and Data Mining*, 2013, p. 158166.
- [16] K. Yu, J. Bi, and V. Tresp, "Active learning via transductive experimental design," in *International Conference on Machine Learning*, 2006, p. 10811088.
- [17] J. Yuan, X. Hou, Y. Xiao, D. Cao, W. Guan, and L. Nie, "Multi-criteria active deep learning for image classification," *Knowledge-Based Systems*, vol. 172, pp. 86–94, 2019.
- [18] D. Gudovskiy, A. Hodgkinson, T. Yamaguchi, and S. Tsukizawa, "Deep active learning for biased datasets via fisher kernel self-supervision," in *Computer Vision and Pattern Recognition*, 2020, pp. 9038–9046.
- [19] G. Hinton and S. Roweis, "Stochastic neighbor embedding," in *International Conference on Neural Information Processing Systems*, 2002, p. 857864.
- [20] E. Elhamifar, G. Sapiro, A. Yang, and S. S. Sarsy, "A convex optimization framework for active learning," in *International Conference on Computer Vision*, 2013, pp. 209–216.
- [21] C. Shui, F. Zhou, C. Gagn, and B. Wang, "Deep active learning: Unified and principled method for query and training," in *International Conference on Artificial Intelligence and Statistics*, 2019, pp. 1–10.
- [22] A. Razavian, J. Sullivan, A. Maki, and S. Carlsson, "A baseline for visual instance retrieval with deep convolutional networks," *ITE Transactions on Media Technology and Applications*, vol. 4, pp. 1–8, 2014.
- [23] A. Babenko and V. Lempitsky, "Aggregating deep convolutional features for image retrieval," *Computer Science*, pp. 1–9, 2015.
- [24] R. Arandjelovi, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN architecture for weakly supervised place recognition," in *Computer Vision and Pattern Recognition*, 2016, pp. 5297–5307.
- [25] L. Zheng, Y. Yang, and Q. Tian, "SIFT Meets CNN: A decade survey of instance retrieval," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 5, pp. 1224–1244, 2018.
- [26] X. Lu, Y. Chen, and X. Li, "Discrete deep hashing with ranking optimization for image retrieval," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 6, pp. 2052–2063, 2020.
- [27] Y. Cao, B. Liu, M. Long, and J. Wang, "HashGAN: Deep learning to hash with pair conditional wasserstein GAN," in *Computer Vision and Pattern Recognition*, 2018, pp. 1287–1296.
- [28] S. Eghbali and L. Tahvildari, "Deep spherical quantization for image search," in *Computer Vision and Pattern Recognition*, 2019, pp. 11 682–11 691.
- [29] B. Klein and L. Wolf, "End-To-End supervised product quantization for image search and retrieval," in *Computer Vision and Pattern Recognition*, 2019, pp. 5036–5045.
- [30] Y. Shen, J. Qin, J. Chen, L. Liu, F. Zhu, and Z. Shen, "Embarrassingly simple binary representation learning," in *International Conference on Computer Vision Workshop*, 2019, pp. 2883–2892.
- [31] F. Shen, Y. Xu, L. Liu, Y. Yang, Z. Huang, and H. T. Shen, "Unsupervised deep hashing with similarity-adaptive and discrete optimization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 3034–3044, 2018.
- [32] E. Yang, T. Liu, C. Deng, W. Liu, and D. Tao, "DistillHash: Unsupervised deep hashing by distilling data pairs," in *Computer Vision and Pattern Recognition*, 2019, pp. 2941–2950.
- [33] Y. Shen, L. Liu, and L. Shao, "Unsupervised binary representation learning with deep variational networks," *International Journal of Computer Vision*, vol. 127, no. 11, pp. 1614–1628, 2019.
- [34] J. Xu, C. Shi, C. Qi, C. Wang, and B. Xiao, "Unsupervised part-based weighting aggregation of deep convolutional features for image retrieval," *AAAI Conference on Artificial Intelligence*, pp. 3671–3679, 2017.
- [35] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [36] M. Tan, R. Pang, and Q. V. Le, "EfficientDet: Scalable and efficient object detection," in *Computer Vision and Pattern Recognition*, 2020, pp. 10 778–10 787.
- [37] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C. Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *European Conference on Computer Vision*, vol. 9905, 2016, pp. 3671–3679.
- [38] Z. Zhao, P. Zheng, S. Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 11, pp. 3212–3232, 2019.
- [39] Y. Yang, Z. Ma, F. Nie, X. Chang, and A. G. Hauptmann, "Multi-class active learning by uncertainty sampling with diversity maximization," *International Journal of Computer Vision*, vol. 113, no. 2, pp. 113–127, 2015.
- [40] O. Sener and S. Savarese, "Active learning for convolutional neural networks: A Core-Set Approach," in *International Conference on Learning Representations*, 2018, pp. 1–12.
- [41] B. Settles, M. Craven, and S. Ray, "Multiple-instance active learning," in *Neural Information Processing Systems*, 2008, pp. 1289–1296.
- [42] A. Freytag, E. Rodner, and J. Denzler, "Selecting influential examples: Active learning with expected model output changes," in *European Conference on Computer Vision*, 2014, pp. 562–577.
- [43] C. Käding, E. Rodner, A. Freytag, and J. Denzler, "Active and continuous exploration with deep neural networks and expected model output changes," *CoRR*, vol. abs/1612.06129, pp. 1–6, 2016.
- [44] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [45] F. Radenovi, G. Tolias, and O. Chum, "CNN image retrieval learns from BoW: Unsupervised fine-tuning with hard examples," in *European Conference on Computer Vision*, 2016, pp. 3–20.
- [46] S. S. Husain and M. Bober, "Improving large-scale image retrieval through robust aggregation of local descriptors," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 9, pp. 1783–1796, 2017.
- [47] Y. Kalantidis, C. Mellina, and S. Osindero, "Cross-dimensional weighting for aggregated deep convolutional features," in *European Conference on Computer Vision Workshops*, 2016, pp. 685–701.
- [48] L. Xie, L. Zheng, J. Wang, A. Yuille, and Q. Tian, "InterActive: Inter-layer activeness propagation," in *Computer Vision and Pattern Recognition*, 2016, pp. 270–279.
- [49] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826.
- [50] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [51] A. Howard, M. Sandler, B. Chen, W. Wang, L. Chen, M. Tan, G. Chu, V. Vasudevan, Y. Zhu, R. Pang, H. Adam, and Q. Le, "Searching for MobileNetV3," in *International Conference on Computer Vision*, 2019, pp. 1314–1324.
- [52] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *Computer Vision and Pattern Recognition*, 2007, pp. 1–8.

- [53] J. Philbin and O. Chum and M. Isard and J. Sivic and A. Zisserman, "Lost in quantization: Improving particular object retrieval in large scale image databases," in *Computer Vision and Pattern Recognition*, 2008, pp. 1–8.
- [54] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," *Computer Science Department, University of Toronto, Tech. Rep.*, vol. 1, pp. 3671–3679, 2009.
- [55] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv:1708.07747*, pp. 1–6, 2017.
- [56] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (VOC) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [57] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.



Wei Huang received the B.Eng. degree in computer application technology from the Henan University of Science and Technology, Luoyang, China, in 2007.

He is currently pursuing the Ph.D. degree in optical engineering with the School of Optical-Electrical and Computer Engineering, University of Shanghai for Science and Technology, Shanghai, China. His current research interests include image processing, deep learning, and computer vision.



Shuzhou Sun received the B.Eng. degree in computer science and technology from Henan Agricultural University, Zhengzhou, China, in 2018.

He is currently pursuing the M.Eng. degree in computer science with the College of Information, Mechanical and Electrical Engineering, Shanghai Normal University, Shanghai, China. His current research interests include computer vision, deep learning, and image processing.



Xiao Lin received the Ph.D. degree in computer science from Shanghai Jiao Tong University, Shanghai, China.

She is currently a Full Professor with the College of Information, Mechanical and Electrical Engineering, Shanghai Normal University (SHNU), Shanghai, China. She is also a Visiting Scholar with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China. Her current research interests include image processing, computer vision, and machine learning.



Ping Li (Member, IEEE) received the Ph.D. degree in computer science and engineering from The Chinese University of Hong Kong, Hong Kong, in 2013.

He is currently an Assistant Professor with the Department of Computing and an Assistant Professor with the School of Design, The Hong Kong Polytechnic University, Hong Kong. He has published many top-tier scholarly research articles, pioneered several new research directions, and made a series of landmark contributions in his areas. He has an

excellent research project reported by the *ACM TechNews*, where only reports the top breakthrough news in computer science worldwide. More importantly, however, many of his research outcomes have strong impacts to research fields, addressing societal needs and contributed tremendously to the people concerned. His current research interests include image/video stylization, colorization, artistic rendering and synthesis, realism in non-photorealistic rendering, computational art, and creative media.



Lei Zhu received the Ph.D. degree in computer science and engineering from The Chinese University of Hong Kong, Shatin, Hong Kong, in 2017.

He is currently working as an Assistant Professor with the ROAS Thrust, HKUST (GZ), and also an affiliated Assistant Professor in ECE with HKUST. Before that, he was a Postdoctoral Researcher with the DAMTP, University of Cambridge, Cambridge, U.K. His current research interests include computer graphics, computer vision, and deep learning.



Jihong Wang received the Ph.D. degree in sport from the Shanghai University of Sport, Shanghai, China, in 2017.

He is currently an Associate Researcher with the Shanghai University of Sport, Shanghai, China. His current research interests include sports engineering, sports health promotion, sports education and training.



C. L. Philip Chen (Fellow, IEEE) received the Ph.D. degree in electrical engineering from Purdue University, West Lafayette, IN, USA, in 1988.

He is currently a Chair Professor and the Dean of the School of Computer Science and Engineering, South China University of Technology, Guangzhou, China. Being a Program Evaluator of the Accreditation Board of Engineering and Technology Education (ABET) in the U.S., for computer engineering, electrical engineering, and software engineering programs, he successfully architects the University

of Macau's Engineering and Computer Science programs receiving accreditations from Washington/Seoul Accord through Hong Kong Institute of Engineers (HKIE), of which is considered as his utmost contribution in engineering/computer science education for Macau as the former Dean of the Faculty of Science and Technology. He is a Fellow of IEEE, AAAS, IAPR, CAA, and HKIE; a member of Academia Europaea (AE), European Academy of Sciences and Arts (EASA), and International Academy of Systems and Cybernetics Science (IASCYS). He received IEEE Norbert Wiener Award in 2018 for his contribution in systems and cybernetics, and machine learnings. He is also a highly cited researcher by Clarivate Analytics in 2018 and 2019.

His current research interests include systems, cybernetics, and computational intelligence. Dr. Chen was a recipient of the 2016 Outstanding Electrical and Computer Engineers Award from his alma mater, Purdue University (in 1988), after he graduated from the University of Michigan at Ann Arbor, Ann Arbor, MI, USA in 1985. He was the IEEE Systems, Man, and Cybernetics Society President from 2012 to 2013, the Editor-in-Chief of the IEEE Transactions on Cybernetics (2020-2021) and the IEEE Transactions on Systems, Man, and Cybernetics: Systems (2014-2019), and currently, an Associate Editor of the IEEE Transactions on Fuzzy Systems. He was the Chair of TC 9.1 Economic and Business Systems of International Federation of Automatic Control from 2015 to 2017, and currently is a Vice President of Chinese Association of Automation (CAA).



Bin Sheng (Member, IEEE) received the Ph.D. degree in computer science and engineering from The Chinese University of Hong Kong, Hong Kong, in 2011.

He is currently a Full Professor with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China. He is an Associate Editor of the IEEE Transactions on Circuits and Systems for Video Technology. His current research interests include virtual reality and computer graphics.