# Universal Domain Adaptation via Compressive Attention Matching

Didi Zhu [*1], Yinchuan Li [*2], Junkun Yuan[1], Zexi Li[1], Kun Kuang [†1], and Chao Wu [†1]

[1]Zhejiang University
[2]Huawei Noah's Ark Lab
[1]{didi_zhu, yuanjk, zexi.li, kunkuang, chao.wu}@zju.edu.cn
[2]{liyinchuan}@huawei.com

## Abstract

*Universal domain adaptation (UniDA) aims to transfer knowledge from the source domain to the target domain without any prior knowledge about the label set. The challenge lies in how to determine whether the target samples belong to common categories. The mainstream methods make judgments based on the sample features, which overemphasizes global information while ignoring the most crucial local objects in the image, resulting in limited accuracy. To address this issue, we propose a Universal Attention Matching (UniAM) framework by exploiting the self-attention mechanism in vision transformer to capture the crucial object information. The proposed framework introduces a novel Compressive Attention Matching (CAM) approach to explore the core information by compressively representing attentions. Furthermore, CAM incorporates a residual-based measurement to determine the sample commonness. By utilizing the measurement, UniAM achieves domain-wise and category-wise Common Feature Alignment (CFA) and Target Class Separation (TCS). Notably, UniAM is the first method utilizing the attention in vision transformer directly to perform classification tasks. Extensive experiments show that UniAM outperforms the current state-of-the-art methods on various benchmark datasets.*

## 1. Introduction

While deep neural networks have achieved remarkable success on visual tasks [12, 25, 19, 73, 51, 69, 70], their performance heavily relies on the assumption of independently and identically distributed (i.i.d.) training and test data [57]. However, this assumption is frequently violated due to the presence of domain shift in real-world scenarios [52, 33, 34, 68, 67, 74, 72]. Unsupervised Domain Adaptation (DA) [1]
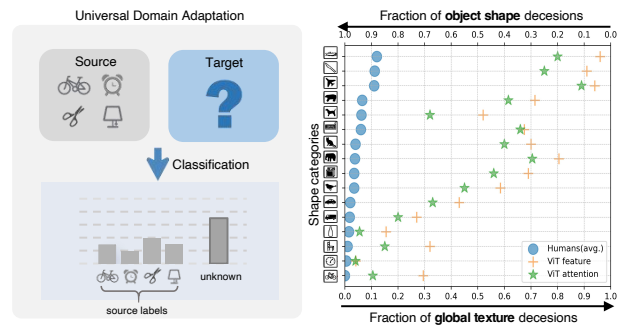


Figure 1: **Left:** Illustration of Universal Domain Adaptation. **Right:** Shape-bias Analysis. Plot shows shape-texture trade off for attention and feature in ViT and humans.

has emerged as a promising solution to address this limitation by adapting models trained on a source domain to perform well on an unlabeled target domain. Nevertheless, most existing DA approaches [14, 56, 49, 40, 39, 38] assume that the label spaces in the source and target domains are identical, which may not always hold in practical scenarios. Partial Domain Adaptation (PDA) [4] and Open Set Domain Adaptation (OSDA) [43] have been proposed to handle cases where the label spaces in one domain include those in the other, but these still rely on prior knowledge on label set, limiting knowledge generalizing from one scenario to others. Universal domain adaptation (UniDA) [66] considers a more practical and challenging scenario where the relationship of label space between source and target domains is completely unknown i.e. with any number of common, source-private and target-private classes.

In UniDA, the primary objective is to develop a model capable of precisely categorizing target samples as one of the common classes or an "unknown" class as shown in Fig. 1 left. Existing UniDA methods aim to design a transferability criteria to detect common and private classes solely based on the discriminability of deep fea-
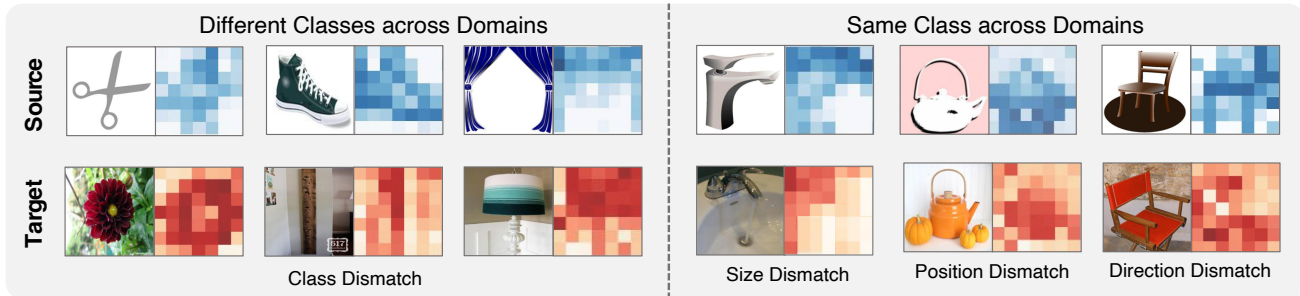
---

Figure 2: **Attention Visualization accross domains.** Attention patterns vary significantly between different classes of images. However, within the same class, attention can also exhibit variations due to differences in object size, position, and orientation. These variations are collectively referred to as *attention mismatch*.

tures [6, 7, 8, 9, 13, 27, 29, 47, 48, 50, 66]. However, over-reliance on deep features can impede model adaptation performance, as they have a strong bias towards global information like texture rather than the essential object information like shape [16, 20], which is considered by humans as the most critical cue for recognition [28]. Fortunately, recent studies have demonstrated that vision transformer (ViT) [24] exhibits a stronger shape bias than Convolutional Neural Network (CNN) [41, 55]. As shown in Fig. 1 right, we confirmed that such strong object shape bias is mainly attributed to the self-attention mechanism, verified in a similar way as [16]. Figure 2 demonstrates the attention vectors of samples across domains. Although we can leverage the attention to focus on more object parts, the *attention mismatch* problem may still exist due to domain shift, which refers to the attention vectors of same-class samples from different domains having some degree of the difference caused by potential variations in object size, orientation, and position across different domains. Attention mismatch can hinder the accurate classification of samples, especially when objects of different classes share similar sizes or positions. For example, in Figure 2, the kettle in the source domain and the flower in the target domain have more similar attention patterns. Therefore, the key challenge in utilizing attention is to effectively explore and leverage the object information embedded in attention while mitigating the negative impact of attention mismatch.

In this paper, we propose a novel Universal Attention Matching (UniAM) framework to address the UniDA problem by leveraging both the feature and attention information in a complementary way. Specifically, UniAM introduces a Compressive Attention Matching (CAM) approach to solve the attention mismatch problem implicitly by sparsely representing target attentions using source attention prototypes. This allows CAM to identify the most relevant attention prototype for each target sample and distinguish irrelevant private labels. Furthermore, a residual-based measurement is proposed in CAM to explicitly distinguish common and

private samples across domains. By integrating attention information with features, we can mitigate the interference caused by domain shift and focus on label shift to some extent. With the guidance of CAM, the UniAM framework achieves domain-wise and category-wise common feature alignment (CFA) and target class separation (TCS). By using an adversarial loss and a source contrastive loss, CFA identifies and aligns the common features across domains, ensuring their consistency and transferability. On the other hand, TCS enhances the compactness of the target clusters, leading to better separation among all target classes. This is accomplished through a target contrastive loss, which encourages samples from the same target class to be closer together and farther apart from samples with other classes.

**Main Contributions**: (1) We propose the UniAM framework that comprehensively considers both attention and feature information, which allows for more accurate identification of common and private samples. (2) We validate the strong object bias of attention in ViT. To the best of our knowledge, we are the first to directly utilize attention in ViT for classification prediction. (3) We implicitly explore object information by sparsely reconstructing attention, enabling better common feature alignment (CFA) and target class separation (TCS). (4) We conduct extensive experiments to show that UniAM can outperform current state-of-the-art approaches.

## 2. Related Works

### 2.1. Universal Domain Adaptation

UniDA [66] does not require prior knowledge of label set relationship. To address this problem, UAN [66] proposes a criterion based on entropy and domain similarity to quantify sample transferability. CMU [13] follows this paradigm to detect open classes by setting the mean of three uncertain scores including entropy, consistency and confidence as a new measurement. Afterward, [27] proposes a real-time adaptive source-free UniDA method. In [47] and

[29], clustering is developed to solve this problem. [31]. OVANet [48] employs a One-vs-All classifier for each class and decides known or unknown by using the output. Recent works have shifted their focus towards finding mutually nearest neighbor samples of target samples [7, 8, 9, **?**] or constructing relationships between target samples and source prototypes [26, 6].

## 2.2. Vision Transformer

Inspired by the success of Transformer [58, 23] in the NLP field, many researchers have attempted to exploit it for solving computer vision tasks. One of the most pioneering works is Vision Transformer (ViT)[24], which decomposes input images into a sequence of fixed-size patches. Different from CNNs that rely on image-specific inductive bias, ViT takes the advantage of large-scale pre-training data and global context modeling on the entire images. Due to the outstanding performance of ViT, many approaches have been proposed based on it [54, 36, 60, 21, 37], such as Touvron et al. [54] propose DeiT, which introduces a distillation strategy specific to transformers to reduce computational costs. In general, ViT and its variants have achieved excellent results on many computer vision tasks, such as object detection [5, 76, 60], image segmentation [75, 61], and video understanding [17, 42], etc.

Recently, ViT has been adopted for the DA task in several works. TVT [65] proposes an transferable adaptation module to capture discriminative features and achieve domain alignment. SSRT [53] formulates a comprehensive framework, which pairs a transformer backbone with a safe self-refinement strategy to navigate challenges associated with large domain gaps effectively. CDTrans [64] designs a triple-branch framework to apply self-attention and cross-attention for source-target domain feature alignment. Differently, we focus in this paper on investigating attention mechanism's superior discriminability across different classes on universal domain adaptation.

## 2.3. Sparse Representation Classification

Sparse Representation Classification (SRC) [62] and Collaborative Representation Classification (CRC)[71], along with their numerous extensions [35, 63, 11, 10, 18], have been extensively investigated in the field of face recognition using single images and videos. These methods have demonstrated promising performance in the presence of occlusions and variations in illumination. By modeling the test data in terms of a sparse linear combination of a dictionary, SRC can capture non-linear relationships between features. Our UniAM is inspired by them but uses a novel measurement instead of a sparsity concentration index.

# 3. Problem Formulation and Preliminary

## 3.1. Problem Formulation

Denoting $\mathbb{X}$, $\mathbb{Y}$, $\mathbb{Z}$ as the input space, label space and latent space, respectively. Elements of $\mathbb{X}$, $\mathbb{Y}$, $\mathbb{Z}$ are noted as $\boldsymbol{x}$, $y$ and $\boldsymbol{z}$. Let $P_s$ and $P_t$ be the source distribution and target distribution, respectively. We are given a labeled source domain $\mathbb{D}_s = \{\boldsymbol{x}_i, y_i)\}_{i=1}^m$ and an unlabeled target domain $\mathbb{D}_t = \{\boldsymbol{x}_i\}_{i=1}^n$ are respectively sampled from $P_s$ and $P_t$, where $m$ and $n$ denote the number of samples of source and target domains, respectively. Denote $\mathbb{L}_s$ and $\mathbb{L}_t$ as the label sets of the source and target domains, respectively. Let $\mathbb{L} = \mathbb{L}_s \cap \mathbb{L}_t$ be the common label set shared by both domains, while $\overline{\mathbb{L}}_s = \mathbb{L}_s \backslash \mathbb{L}$ and $\overline{\mathbb{L}}_t = \mathbb{L}_t \backslash \mathbb{L}$ be the label sets private to source and target domains, respectively. Denote $M = |\mathbb{L}_s|$ as the number of source labels. Universal domain adaptation aims to predict labels of target data in $\mathbb{L}$ while rejecting the target data in $\overline{\mathbb{L}}^t$ based on $\mathbb{D}_s$ and $\mathbb{D}_t$.

Our overall architecture consists of a ViT-based feature extractor, an adversarial domain classifier, and a label classifier. Suppose the function for learning embedding features is $G_f : \mathbb{X} \to \mathbb{Z} \in \mathbb{R}^{d_z}$ where $d_z$ is the length of each feature vector, the discrimination function of the label classifier is $G_c : \mathbb{Z} \to \mathbb{Y} \in \mathbb{R}^M$, and the function of the domain classifier is $G_d : \mathbb{Z} \to \mathbb{R}^1$.

## 3.2. Preliminary

To start with, we provide an overview of the self-attention mechanism used in ViT. First, the input image $\boldsymbol{x}$ is divided into $N$ fixed-size patches, which are linearly embedded into a sequence of vectors. Next, a special token called the class token is prepended to the sequence of image patches for classification. The resulting sequence of length $N + 1$ is then projected into three matrices: queries $\boldsymbol{Q} \in \mathbb{R}^{(N+1) \times d_k}$, keys $\boldsymbol{K} \in \mathbb{R}^{(N+1) \times d_k}$ and values $\boldsymbol{V} \in \mathbb{R}^{(N+1) \times d_v}$ with $d_k$ and $d_v$ being the length of each query and value vector, respectively. Then, $\boldsymbol{Q}$ and $\boldsymbol{K}$ are passed to the self-attention layer to compute the patch-to-patch similarity matrix $\boldsymbol{A}^{(N+1) \times (N+1)}$, which is given by

$$\boldsymbol{A} = \frac{\boldsymbol{Q}\boldsymbol{K}^\top}{\sqrt{d_k}}, \tag{1}$$

For ease of further processing, we flatten $\boldsymbol{A}$ into a vector $\boldsymbol{a} \in \mathbb{R}^{(N+1)^2 \times 1}$. It is worth noting that multiple attention heads are utilized in the self-attention mechanism. Each head outputs a separate attention, and the final attention is obtained by concatenating the vectors from all heads. As a result, the dimensionality of $\boldsymbol{a} \in \mathbb{R}^{d_a \times 1}$, where $d_a = N_H \times (N + 1)^2$ and $N_H$ is the number of attention heads. The utilization of multiple heads allows the model to jointly attend to information from different feature subspaces at different positions.
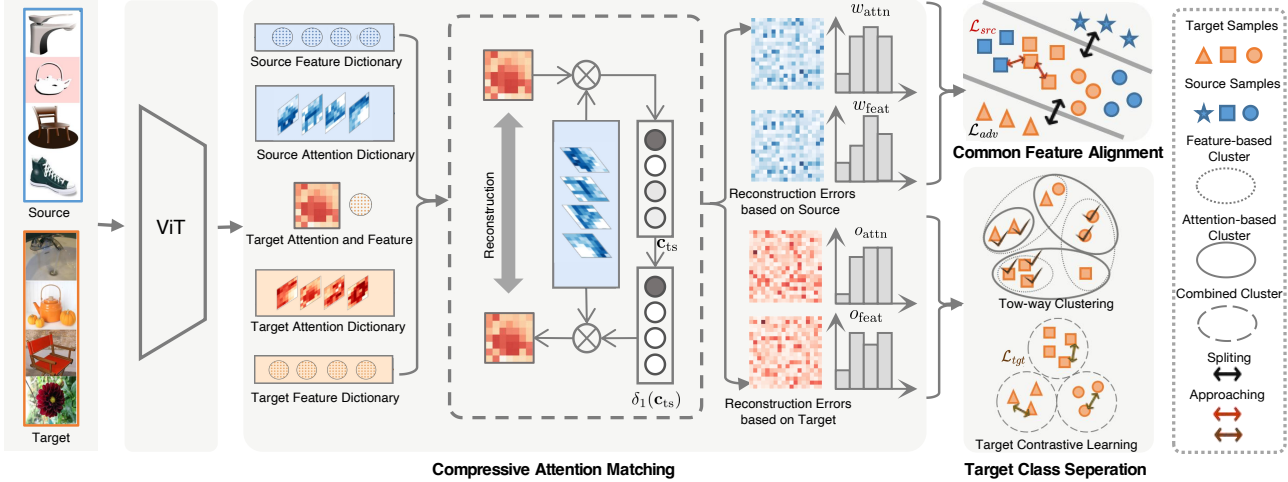
Figure 3: **Illustration of the proposed UniAM framework.** The framework consists of three integral components: Compressive Attention Matching (CAM), Common Feature Alignment (CFA) and Target Class Separation (TCS). At its core, CAM reconstructs all target attentions and features based on the source dictionary (with feature reconstruction omitted in the figure for simplicity), and attention and feature commonness scores $w_{attn}$ and $w_{feat}$ are computed from residual vectors. Then, domain- and category-wise CFA is achieved by minimizing $\mathcal{L}_{adv}$ and $\mathcal{L}_{src}$ guided by $w_{attn}$ and $w_{feat}$. Similarly, $o_{attn}$ and $o_{feat}$ are obtained by reconstructing target attentions and features based on the target dictionary in CAM. TCS performs two-way clustering from both the attention and feature views and minimizes $\mathcal{L}_{tgt}$ to achieve effective separation of target classes.

Once the attention vector $a$ is available, the corresponding $k$-th attention prototype $p_k$ is calculated by averaging all attention vectors of samples in class $k$, which will be used in the subsequent matching process.

## 4. Proposed Methodology

### 4.1. Compressive Attention Matching

Since the attention mismatch problem exists due to domain shift mentioned in Section 1, how to effectively utilize the core object information and avoid interference from redundant information poses a challenge in applying attention to UniDA. To address this challenge, compressive attention matching (CAM) is proposed to capture the most informative object structures by sparsely representing target attentions. Define the attention dictionary in CAM as the collection of source attention prototypes for efficient matching, i.e., $\boldsymbol{P}_s = [\boldsymbol{p}_1^s, \boldsymbol{p}_2^s, \cdots, \boldsymbol{p}_M^s] \in \mathbb{R}^{d_a \times M}$. Definition 1 gives the definition of CAM.

**Definition 1** (**Compressive Attention Matching**). *Given an attention vector $\boldsymbol{a}_t \in \mathbb{R}^{d_a \times 1}$ of the target sample $\boldsymbol{x}_t$ and a source attention dictionary $\boldsymbol{P}_s$, Compressive Attention Matching aims to match $\boldsymbol{a}_t$ with one prototype in $\boldsymbol{P}_s$ to determine its commonness, which is achieved by assuming that $\boldsymbol{a}_t$ can be approximated by a linear combination of $\boldsymbol{P}_s$:*

$$\boldsymbol{a}_t = \boldsymbol{P}_s \boldsymbol{c}_{ts}, \qquad (2)$$

*where the coefficient vector $\boldsymbol{c}_{ts} \in \mathbb{R}^{M \times 1}$ satisfies a **sparsity constraint** in order to achieve a compressive representation. Based on $\boldsymbol{c}_{ts}$, $\boldsymbol{x}_t$ is regarded as belonging to common classes from an attention perspective when the following inequality is satisfied:*

$$w_{attn}(\boldsymbol{x}_t) < \delta.$$

$w_{attn}(\cdot)$ *indicates a measurement to evaluate the commonness of $\boldsymbol{x}_t$ which is defined later and $\delta$ is a threshold.*

**Why Compressive Attention Matching is desirable?** By enforcing sparsity on the coefficients in CAM, we can obtain a compressive representation of the attention vectors, which facilitates the extraction and utilization of low-dimensional structures embedded in high-dimensional attention vectors. In the context of UniDA, this compressive representation enables us to identify the most relevant attention prototype for each target sample and distinguish irrelevant private labels, which is crucial for achieving effective common and private class detection. Therefore, CAM with sparse coefficients plays a vital role in solving UniDA.

To solve Eq. 2 in CAM, the coefficient vector $\boldsymbol{c}_{ts}$ is estimated by:

$$\min_{\boldsymbol{c}_{ts}} \|\boldsymbol{a}_t - \boldsymbol{P}_s \boldsymbol{c}_{ts}\|_2^2 + \rho \|\boldsymbol{c}_{ts}\|_1, \qquad (3)$$

where $\|\cdot\|_p$ denotes $\ell_p$-norm. The $\ell_1$-minimization term in Eq. 3 yields a sparse solution, which enforces that $\boldsymbol{c}_{ts}$ has only a small number of non-zero coefficients.

Then we can compute the class reconstruction error vector $\boldsymbol{r}_{ts} \in \mathbb{R}^M$ for each target sample using the sparse matrix $\boldsymbol{c}_{ts}$. The $k$-th entry of $\boldsymbol{r}_{ts}$ can be represented:

$$\boldsymbol{r}_{ts}(k) = \|\boldsymbol{a}_t - \boldsymbol{P}_s \delta_k(\boldsymbol{c}_{ts})\|_2, \quad k = 1, \ldots, M, \quad (4)$$

where $\delta_k(\boldsymbol{c}_{ts})$ is a one-hot vector with the $k$-th entry in $\boldsymbol{c}_{ts}$ being non-zero while setting all other entries to zero. If $\boldsymbol{x}_t$ corresponds to a common class $k$, then the reconstruction error corresponding to class $k$. $\boldsymbol{r}_{ts}(k)$ should be much lower than that corresponding to the other classes. Conversely, if $\boldsymbol{x}_t$ belongs to a private class, the difference between elements of the entire reconstruction error vector $\boldsymbol{r}_{ts}$ should be relatively small, without a significant difference between the errors corresponding to different classes.

As a result, the reconstruction error vector $\boldsymbol{r}_{ts}$ is a crucial component in CAM. It serves as the foundation for the design of the measurement $w_{\text{attn}}(\cdot)$ in Definition 1 called Attention Commonness Degree (ACD), defined as belows:

**Definition 2 (Attention Commonness Degree).** *Given the residual vector $\boldsymbol{r}_{ts}$ of $\boldsymbol{x}_t$, the ACD is defined as the difference between the average of non-matched errors and matched errors:*

$$w_{attn}(\boldsymbol{x}_t) = non\text{-}match(\boldsymbol{r}_{ts}) - match(\boldsymbol{r}_{ts}), \quad (5)$$

*where $match(\boldsymbol{r}_{ts}) = \boldsymbol{r}_{ts}(\hat{y})$, $\hat{y} = \arg\min_k \boldsymbol{r}_{ts}(k)$ and non-match$(\boldsymbol{r}_{ts})$ is the average of reconstruction errors excepting $\hat{y}$.*

**Remark 1.** *ACD measures the degree of commonness for a target sample $\boldsymbol{x}_t$, which represents the probability of belonging to common classes. A higher ACD value indicates a larger difference between non-matched and matched errors, suggesting the presence of an attention prototype similar to $\boldsymbol{x}_t$, and consequently, a higher degree of sample commonness. Conversely, a smaller ACD value implies a similar reconstruction error between $\boldsymbol{x}_t$ and all source prototypes, indicating a lower degree of sample commonness and a higher degree of privateness.*

To complement the attention information, we retain features that reflect global information. The target feature $\boldsymbol{z}_t$ can be also represented by the linear span of source feature prototypes $\boldsymbol{Q}_s = [\boldsymbol{q}_1^s, \boldsymbol{q}_2^s, \cdots, \boldsymbol{q}_M^s]$, i.e., $\boldsymbol{z}_t = \boldsymbol{Q}_s \boldsymbol{c}_{ts}$. The corresponding residual vector $\boldsymbol{r}'_{ts}(k)$ is computed based on $\boldsymbol{c}_{ts}$. The Feature Commonness Degree (FCD) can be defined as $w_{\text{feat}}(\boldsymbol{x}_t) = \text{non-match}(\boldsymbol{r}'_{ts}) - \text{match}(\boldsymbol{r}'_{ts})$..

It is worth noting that by replacing $P_s$ in Definition 2 with the target dictionary $P_t$, we can obtain compressive representations of target attentions towards $P_t$. This leads to a score similar to that in Definition 2, denoted as $o_{\text{attn}}$. The same goes for $o_{\text{feat}}$. These scores can facilitate determining the probability that two target samples belong to the same class, more details will be provided in Section 4.3.

In summary, both attention and feature characteristics are important factors that affect the perception of similarity between different categories. Attention captures the structural properties of objects, while feature captures the appearance properties of the global images. Therefore, we can achieve a more comprehensive and accurate private class detection model that takes into account both object information and global information.

### 4.2. Common Feature Alignment

To identify and align the common class features across domains, we propose a domain-wise and category-wise Common Feature Alignment (CFA) technique, which considers both attention and feature information.

**Domain-wise Alignment.** To achieve domain-wise alignment, we first propose a residual-based transferability score $d_t$ measuring the probability that the target sample belongs to the common classes, which can be summarized as:

$$w_t = \lambda w_{\text{attn}} + (1 - \lambda) w_{\text{feat}}, \quad (6)$$

where $\lambda$ is a hyperparameter balancing their contribution. Meanwhile, to measure the probability that the source sample $\boldsymbol{x}_s$ with label $j$ belongs to the common label set, we compute $w_j^s$ with the sum of all target samples' attention and feature reconstruction errors respectively, i.e.

$$w_s^j = \lambda \sigma(\boldsymbol{r}_{ts})_j + (1 - \lambda)\sigma(\boldsymbol{r}'_{ts})_j \quad (7)$$

where $\boldsymbol{r}_{ts}^i$ indicates the reconstruction error of the $i$-th target sample. The operator $\sigma(\boldsymbol{r}) = \frac{\boldsymbol{r} - min(\boldsymbol{r})}{max(\boldsymbol{r}) - min(\boldsymbol{r})}$ refers to the normalization sum of all target attention or feature residual vectors. A larger value of $w_j^s$ indicates a higher probability that the source label $j$ belongs to the common label set, while lower values suggest that it is more likely to be a source private label. It is worth noting that source samples with the same label are assigned the same weight.

Based on the above two weights, we can derive a domain-wise adversarial loss that aligns the common classes across domains as follows:

$$\mathcal{L}_{\text{adv}} = \mathbb{E}_{\boldsymbol{x}_s \in \mathbb{D}_S}\left[w_s \cdot \log\left(1 - G_d\left(\boldsymbol{z}_s\right)\right)\right] \\ + \mathbb{E}_{\boldsymbol{x}_t \in \mathbb{D}_T}\left[w_t \cdot \log\left(G_d\left(\boldsymbol{z}_t\right)\right)\right], \quad (8)$$

In addition, to avoid being interfered with the knowledge of source private samples, we employ an indicator as the weight for the weighted cross-entropy loss $\mathcal{L}_{\text{cls}}$ for the source domain, as shown below:

$$\mathcal{L}_{\text{cls}} = -\mathbb{E}_{(\boldsymbol{x}_s, y_s) \in \mathbb{D}_S} \mathbb{1}_{w_s^y \geqslant \alpha} l_{ce}(y, G_c(G_f(\boldsymbol{x}_s))), \quad (9)$$

where $l_{ce}$ is the standard cross-entropy loss.

**Category-wise Alignment.** To enhance the source discriminability and align the common features from a category-wise perspective across domains, we propose a contrastive

common feature alignment method. In order to quantify the likelihood that the source sample $x_i^s$ and the sample $x_j$ belong to the same category $y_i$, we design a category-wise target score $w_{i,j}$. If $x_j$ is a source sample, we use its ground truth label $y_j$ to determine if it's a positive or negative example, i.e., $w_{i,j} = \begin{cases} 1, & \text{if } y_j = y_i \\ 0, & \text{if } y_j \neq y_i \end{cases}$. If $x_j$ is a target sample, we estimate the probability that it belongs to $y_i$ based on the residual vectors $r_{ts}$ and $r'_{ts}$. These two vectors can be seen as two vanilla prediction probability vectors and the corresponding pseudo-labels $\hat{y}_j$ and $\hat{y}'_j$ can be obtained by argmin operation. To give a more reliable estimation, the soft label of $x_j$ is determined by these two pseudo-labels together. Specifically, $w_{i,j}$ is set to 1 when both $\hat{y}_j$ and $\hat{y}'_j$ are equal to $y_i$ and set to 0 when both of them are not equal to $y_i$. The soft label is computed as below when only one of the predictions is $y_i$:

$$w_{i,j} = \begin{cases} \lambda \cdot w_{\text{attn}}/w_t, & \text{if } \hat{y}'_j \neq \hat{y}_j = y_i, \\ (1-\lambda) \cdot w_{\text{feat}}/w_t, & \text{if } \hat{y}_j \neq \hat{y}'_j = y_i. \end{cases} \quad (10)$$

Thus the category-wise common feature alignment can be improved by minimizing the source contrastive loss $\mathcal{L}_{\text{src}}$:

$$\mathcal{L}_{\text{src}} = -\mathbb{E}_{x_i \in \mathbb{D}_S, x_j \in \mathbb{D}_{S \cup T}} w_{i,j} l(z_i, z_j) \quad (11)$$

with

$$l(z_i, z_j) = \frac{\exp(z_i z_j/\tau)}{\sum_{k=1}^{m+n} \exp(z_i z_k/\tau)}.$$

### 4.3. Target Class Seperation

To better distinguish the common and private classes in the target domain, we propose a Target Class Separation (TCS) technique. As mentioned in Section 4.1, we can construct a CAM problem based on the target attention dictionary $P_t$. As the target labels are unknown, $P_t = [p_1^t, p_2^t, \cdots, p_K^t] \in \mathbb{R}^{d_a \times K}$ is initialized by performing traditional K-means algorithm on target attentions $\{a_t^i\}_{i=1}^n$, where $K$ is pre-defined. In the subsequent attention clustering process, we calculate the residual vector $r_{tt} \in \mathbb{R}^K$ and use it as a metric for measuring the distance between samples, which allows for dynamically update $P_t$. Iteratively refining $P_t$ makes it more reliable and discriminative. Meanwhile, the feature clustering based on the target feature dictionary $Q_t = [q_1^t, q_2^t, \cdots, q_K^t] \in \mathbb{R}^{d_z \times K}$ is also performed. After the two-way clustering, each target sample $x_i^t$ is assigned two cluster indexes $\hat{c}_i$ and $\hat{c}'_i$ from attention and feature view like Fig. **??** depicted. The final soft pseudo label $o_{c,i}$ determining whether $x_i$ belong to the $c$-th cluster is obtained based on these two cluster indexes, similar to $w_{i,j}$. Based on $o_{i,j}$, the target contrastive loss $\mathcal{L}_{\text{tgt}}$ is computed as follows:

$$\mathcal{L}_{\text{tgt}} = -\mathbb{E}_{x_i \in \mathbb{D}_T, x_j \in \mathbb{D}_T} o_{i,j} l(z_i, z_j), \quad (12)$$

with

$$l(z_i, z_j) = \frac{\exp(z_i z_j/\tau)}{\sum_{k=1}^n \exp(z_i z_k/\tau)}.$$

where $o_{i,j} = o_{c,i} \cdot o_{c,j}$ is the probability weight determining whether $x_i$ and $x_j$ belong to the same cluster $c$. By minimizing $\mathcal{L}_{\text{tgt}}$, we can enhance the compactness of target clusters making a better separation among target classes.

### 4.4. Overall Framework

Overall, our framework is jointly optimized with four terms, i.e., cross-entropy loss $\mathcal{L}_{\text{cls}}$, adversarial loss $\mathcal{L}_{\text{adv}}$, source and target contrastive loss $\mathcal{L}_{\text{src}}$ and $\mathcal{L}_{\text{tgt}}$ as shown in Fig. 3,

$$\max_{G_d} \min_{G_f, G_c} \mathcal{L}_{\text{cls}} + \eta_1 \mathcal{L}_{\text{src}} + \eta_2 \mathcal{L}_{\text{tgt}} - \mathcal{L}_{\text{adv}}, \quad (13)$$

where $\eta_1$ and $\eta_2$ are set as 0.5 to balance each loss component. In the testing phase, given each input target sample $x_t$, we compute $w_t$ in (6). For those samples that satisfy $w_t < \beta$ are assigned with the predicted source class, where $\beta$ is a validated threshold. Otherwise, the samples are marked as unknown.

## 5. Experiment Results

### 5.1. Experimental Setup

**Datasets.** We perform experiments on **Office-31** [46], **Office-Home** [59], **VisDA2017** [45] and **DomainNet** [44] datasets. **Office-31** consists of three domains: Amazon (A), DSLR (D) and Webcam (W). Each domain contains 31 categories. **Office-Home** is a dataset made up of 65 different categories from four domains: Artistic (Ar), Clipart (Cl), Product (Pr) and Real-world images (Rw). **VisDA2017** is a dataset with a single source and target domain testing the ability to perform transfer learning from synthetic images to natural images. The dataset has 12 categories in each domain. We conduct experiments on three subsets from it, i.e., Painting (P), Real (R), and Sketch (S). For a fair comparison, we follow the same dataset split as [66] for the first three dataset and [13] for the last dataset.

**Evaluation Protocols.** We evaluate all methods using H-score [13]. H-score is the harmonic mean of the accuracy of common classes and the accuracy of the "unknown" classes, which can make a trade-off between the accuracy of known and unknown classes.

**Implementation Details** The method is implemented in Pytorch using a ViT-base model with $16 \times 16$ input patch size (or ViT-B/16) [24], pretrained on ImageNet [12]), as the backbone feature extractor. The transformer encoder of ViT-B/16 contains a total of 12 Transformer layers. The label classifier consists of a fully connected network with BatchNorm [22]. The domain discriminator is a three-layer MLP with ReLU activations. We train all models using

Table 1: **H-score (%) on Office-31 and DomainNet**

| Method | | **Office-31** | | | | | | | **DomainNet** | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | A2W | D2W | W2D | A2D | D2A | W2A | Avg | P2R | R2P | P2S | S2P | R2S | S2R | Avg |
| ResNet [19] | ResNet50 | 47.92 | 54.94 | 55.60 | 49.78 | 48.48 | 48.96 | 50.94 | 30.06 | 28.34 | 26.95 | 26.95 | 26.89 | 29.74 | 28.15 |
| DANN [15] | | 48.82 | 52.73 | 54.87 | 50.18 | 47.69 | 49.33 | 50.60 | 31.18 | 29.33 | 27.84 | 27.84 | 27.77 | 30.84 | 29.13 |
| OSBP [50] | | 50.23 | 55.53 | 57.20 | 51.14 | 49.75 | 50.16 | 52.34 | 33.60 | 33.03 | 30.55 | 30.53 | 30.61 | 33.65 | 32.00 |
| UAN [66] | | 58.61 | 70.62 | 71.42 | 59.68 | 60.11 | 60.34 | 63.46 | 41.85 | 43.59 | 39.06 | 38.95 | 38.73 | 43.69 | 40.98 |
| CMU [13] | | 67.33 | 79.32 | 80.42 | 68.11 | 71.42 | 72.23 | 50.78 | 52.16 | 45.12 | 44.82 | 45.64 | 50.97 | 48.25 | 73.14 |
| DCC [29] | | 78.54 | 79.29 | 88.58 | 88.50 | 70.18 | 75.87 | 80.16 | 56.90 | 50.25 | 43.66 | 44.92 | 43.31 | 56.15 | 49.20 |
| OVANet [48] | | 79.45 | 95.43 | 94.35 | 85.67 | 80.43 | 84.23 | 86.59 | 56.0 | 51.7 | 47.1 | 47.4 | 44.9 | 57.2 | 50.7 |
| UniOT [6] | | 89.16 | 98.93 | 96.87 | 86.35 | 89.85 | 88.08 | 91.54 | 59.30 | 47.79 | 51.79 | 46.81 | 48.32 | 58.25 | 52.04 |
| OVANet* | ViT | 87.75 | 93.14 | 85.72 | 82.96 | **92.67** | 91.25 | 88.92 | 71.24 | 61.14 | 51.28 | 55.30 | 47.51 | 66.48 | 58.83 |
| UniOT* | | 96.35 | 99.13 | 99.43 | 88.40 | 89.67 | **93.81** | 94.47 | 72.40 | 59.47 | 49.30 | 56.86 | 47.38 | 69.43 | 59.14 |
| Ours | | **95.46** | **99.62** | **99.81** | **95.28** | 92.35 | 93.23 | **95.95** | **73.87** | **60.89** | **52.31** | **59.98** | **51.41** | **70.68** | **61.52** |

Table 2: **H-score (%) on Office-Home and VisDA2017**

| Method | | **Office-Home** | | | | | | | | | | | | | **VisDA** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Ar2Cl | Ar2Pr | Ar2Rw | Cl2Ar | Cl2Pr | Cl2Rw | Pr2Ar | Pr2Cl | Pr2Rw | Rw2Ar | Rw2Cl | Rw2Pr | Avg | S2R |
| ResNet [19] | ResNet50 | 44.65 | 48.04 | 50.13 | 46.64 | 46.91 | 48.96 | 47.47 | 43.17 | 50.23 | 48.45 | 44.76 | 48.43 | 47.32 | 25.44 |
| DANN [15] | | 42.36 | 48.02 | 48.87 | 45.48 | 46.47 | 48.37 | 45.75 | 42.55 | 48.70 | 47.61 | 42.67 | 47.40 | 46.19 | 25.65 |
| OSBP [50] | | 39.59 | 45.09 | 46.17 | 45.70 | 45.24 | 46.75 | 45.26 | 40.54 | 45.75 | 45.08 | 41.64 | 46.90 | 44.48 | 27.31 |
| UAN [66] | | 51.64 | 51.70 | 54.30 | 61.74 | 57.63 | 61.86 | 50.38 | 47.62 | 61.46 | 62.87 | 52.61 | 65.19 | 56.58 | 30.47 |
| CMU [13] | | 56.02 | 56.93 | 59.15 | 66.95 | 64.27 | 67.82 | 54.72 | 51.09 | 66.39 | 68.24 | 57.89 | 69.73 | 61.60 | 34.64 |
| DCC [29] | | 57.97 | 54.05 | 58.01 | 74.64 | 70.62 | 77.52 | 64.34 | **73.60** | 74.94 | 80.96 | **75.12** | 80.38 | 70.18 | 43.02 |
| OVANet [48] | | 62.81 | 75.54 | 78.59 | 70.72 | 68.78 | 75.03 | 71.27 | 58.64 | 80.52 | 76.09 | 64.13 | 78.91 | 71.75 | 53.10 |
| UniOT [6] | | 67.27 | 80.54 | 86.03 | 73.51 | 77.33 | 84.28 | 75.54 | 63.33 | 85.99 | 77.77 | 65.37 | 81.92 | 76.57 | 57.32 |
| OVANet* | ViT | 58.09 | 86.06 | 89.38 | 81.86 | 81.03 | 86.22 | 84.49 | 57.06 | 88.54 | 83.67 | 57.32 | 86.67 | 77.45 | 56.98 |
| UniOT* | | 63.77 | **88.19** | 90.23 | 74.99 | 81.02 | 84.55 | 78.91 | 61.29 | 87.60 | 82.38 | 63.70 | 88.30 | 78.40 | 63.25 |
| Ours | | **72.04** | 87.07 | **90.67** | 80.30 | **82.39** | 79.81 | **85.02** | 68.35 | **88.98** | **85.44** | 72.11 | 86.12 | **81.68** | **65.18** |

a minibatch Stochastic Gradient Descent (SGD) optimizer with a momentum of 0.9 and a weight decay of $5 \times 10^{-4}$. The learning rate decays by a factor of $(1+\alpha i/N)^{-\beta}$, where $i$ and $N$ respectively denote the current iteration and the global iteration. The batch size is set to 36. We initialize the initial learning rate to 0.01 for Office-31 and Office-Home, while set 0.001 for VisDA2017 and DomainNet. For the regularization hyperparameters, we set $\gamma = 100$ and $\lambda = 0.3$ for all dataset. For the decision threshold, we set $\alpha = 0.85$ and $\beta = 1.0$ for all dataset in UniDA and OSDA. In PDA, we set $\alpha = 0.8$ for all dataset except $\alpha = 0.85$ in the Office-31 W2A task. For the pre-defined number of target prototypes, a larger size of the target domain indicates a larger $K$. Therefore, we empirically set $K = 50$ for Office-31, $K = 150$ for Office-Home, $K = 500$ for VisDA, $K = 1000$ for DomainNet.

### 5.2. Comparison Baselines

Follow the previous existing works [13], we compare our method with (1) ResNet [19], (2) close-set domain adapta-

tion: DANN [15], (3) partial domain adaptation: PADA [4], ETN [3], BA³US [30] (4) open set domain adaptation: OSBP [50], STA [32], ROS [2]. (5)universal domain adaptation: UAN [66], CMU [13], DANCE [47], DCC [29], OVANet [48], UniOT [6], GATE [9]. We use some results from [9]. In all experiments, we assume that none of the UniDA methods have prior knowledge of category shift, while baselines tailored for each setting consider this prior.

### 5.3. Comparison Results

The experimental results for the Office-31, Office-Home, VisDA2017, and DomainNet datasets are presented in Table 1 and Table 2, which demonstrate that our proposed UniAM framework outperforms the state-of-the-art approaches in all benchmarks, as evaluated by the H-score metric. Additionally, to ensure a fair comparison, we conducted experiments by replacing the backbone of OVANet and UniOT with ViT, marked as ⋆. The proposed method consistently surpasses these ViT-based methods by a large margin. This indicates that our approach does not solely
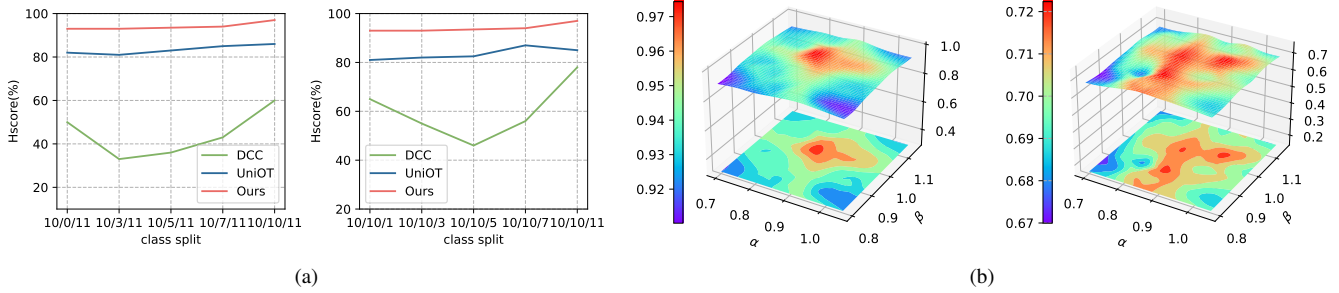
(a)

(b)

Figure 4: (a) **Effectiveness on different label set relationships.** Our method consistently outperforms all comparative approaches across various settings of $\overline{\mathbb{L}}^t$. (b) **Effectiveness of varying decision threshold $\alpha$ and $\beta$.** Even when both $\alpha$ and $\beta$ undergo significant fluctuations, the performance of our method doesn't decline sharply.



Figure 5: **Qualitative analysis of feature and attention.**

Table 3: **Evaluation of the effectiveness of UniAM.**

| Method | A2W | D2W | W2D | A2D | D2A | W2A | Avg |
|--------|-----|-----|-----|-----|-----|-----|-----|
| w/ $\mathcal{L}_{adv}$ | 92.12 | 93.37 | 99.49 | 93.54 | 88.32 | 92.77 | 93.27 |
| w/ $\mathcal{L}_{src}$ | 89.61 | 98.58 | 99.57 | 91.65 | 91.35 | 92.73 | 93.91 |
| w/ $\mathcal{L}_{tgt}$ | 93.26 | 98.27 | 99.78 | 92.10 | 89.21 | 89.97 | 93.76 |
| w/o $\mathcal{L}_{adv}$ | 94.69 | 98.10 | 99.78 | **96.12** | 92.05 | 92.35 | 95.51 |
| w/o $\mathcal{L}_{src}$ | 93.78 | 98.43 | 99.78 | 94.69 | 91.94 | 93.47 | 95.34 |
| w/o $\mathcal{L}_{tgt}$ | 90.76 | 98.20 | 99.78 | 92.41 | 91.69 | 93.13 | 94.32 |
| w/o $w_{attn}$ | 94.76 | 97.54 | 97.26 | 94.31 | 91.22 | 90.69 | 94.30 |
| Ours | **95.46** | **99.62** | **99.81** | 95.28 | **92.35** | **93.23** | **95.95** |

rely on using ViT as the backbone, but rather it fully exploits the advantage of the attention mechanism in ViT for UniDA tasks to achieve such superior performance.

### 5.4. Analysis on Different Label Set Relationships

**Varying size of target private label set $\overline{\mathbb{L}}^t$.** To explore the performance of our method under different class splitting settings with OVANet, UniOT and UniOT⋆, we fix $\mathbb{L}^s$, $\mathbb{L}$ and change $\overline{\mathbb{L}}^t$ on task A→W in Office-31 dataset. As shown in Fig. 4 (a) left, our method consistently outperforms all comparison methods under different $\overline{\mathbb{L}}^t$, proving that our method is effective and robust for different $\overline{\mathbb{L}}^t$. As $\overline{\mathbb{L}}^t$ increases, meaning there are many open classes, our method outperforms other methods by a large margin, demonstrating that our method is superior in detecting open classes.

**Varying size of common label set $\mathbb{L}$** We fix $\mathbb{L}^s$ and $\mathbb{L}^t$ and varying $\mathbb{L}$ on task A→W in Office-31 dataset. We let $\overline{\mathbb{L}}^s$, $\overline{\mathbb{L}}^t$ to keep 10 and 11 and vary $\mathbb{L}$ from 0 to 10 . In particular, all target data should be marked as "unknown" when

the source and target domains do not overlap on label sets. As shown in Fig. 4 (a) right, our method consistently outperforms previous methods on all sizes of $\mathbb{L}$, indicating that our method can detect open classes more effectively.

### 5.5. Analysis on Our Method

**Effectiveness of different losses.** As there are three losses excluding classification loss in our method, we conduct another experiment to verify the effectiveness of each loss and any combination of them on Office-31 dataset. As shown in the first six rows of Table. 3, the results indicate that the use of any single loss function or a combination of any two loss functions can lead to a decrease in performance to some extent, with a performance drop of 2.05%-2.68% observed when using a single loss. These findings demonstrate the effectiveness of our proposed method.

**Effectiveness of attention in ViT.** To demonstrate the indispensable roles of attention in ViT in the proposed transferability criteria, we introduce a variant denoted as w/o $w_{attn}$, which performs sparse reconstruction only on the features. Compared with our method in the last two rows of Table 3, the average performance drop of w/o $w_{attn}$ 1.65%. It indicates that the attention mechanism does play an effec-
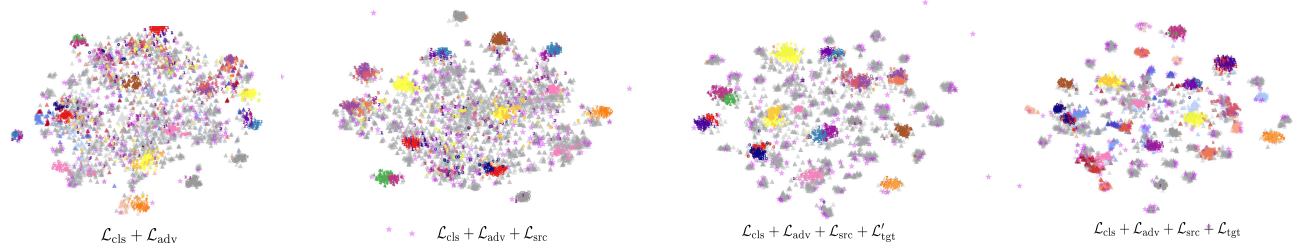
Figure 6: **Feature visualization of target domain with different losses.** $\mathcal{L}'_{tgt}$ and $\mathcal{L}_{tgt}$ refer to using only $w_{feat}$ to calculate the loss weight $w_{i,j}$ and using both $w_{attn}$ and $w_{feat}$ together to calculate $w_{i,j}$. (a-c) $\mathcal{L}src$ increases the distance between common and private categories, treating all target-private samples as one class, while $\mathcal{L}tgt$ enhances discriminability of target private classes. (d) UniAM leverages attention to guide the refinement of target representations.

tive role in attention enhancement on the basis of features during the process of common category detection.

**Qualitative Analysis.** As shown in Fig. 5, the three probability density histograms visualize partially $w_{attn}$, $w_{text}$, and their weighted sum $w_t$ in A→W task on Office. From Fig. 5, it can be observed that using $w_{attn}$ or $w_{text}$ alone can partially distinguish common samples (colored in blue) and private samples (colored in red), but each has its limitations. $w_{attn}$ is prone to confusion at the boundary, while $w_{text}$ has some outliers, such as private samples with extremely high values and common samples with extremely low values. By combining them together, these two limitations can be effectively alleviated. The weighted sum $w_t$ can result in clearer boundaries between private and common samples, and the outliers are reduced.

**Feature visualization.** We use t-SNE to visualize the learned target features for Pr→Rw of Office-Home. As shown in Fig. 6, the gray dots represent private samples, while the non-gray dots represent common samples, and their colors indicate their ground-truth classes. Fig. 6 (a)-(c) shows that $\mathcal{L}_{src}$ increased the distance between common and private categories while all target-private samples are treated as a single class, and $\mathcal{L}_{tgt}$ improved the discriminability of the target private classes. Especially, Fig. 6 (d) validates that UniAM learns a better target representation introducing attention as a guide for attention enhancement can further improve the discriminability in the target domain by bringing same-class samples closer and pushing different-class samples farther away.

**Sensitivity to decision threshold.** We investigate the sensitivity of thresholds $\alpha$ and $\beta$, which are used to determine whether source and target samples belong to common classes respectively. The analysis was done in A→D on Office-31 and Ar→Cl on Office-Home. As depicted in Fig. 4 (b), the H-score demonstrates minimal variance. Specifically, $\alpha$ varies within a reasonable and practical range of [0.7, 1.0], while $\beta$ varies in a range of [0.8, 1.1]. These findings collectively reinforce the idea that our method remains robust to variations in the $\alpha$ and $\beta$ parameters. This robustness is a strong indicator of the method's stability and resilience under varying parameter settings.

# 6. Conclusions

In this work, we introduced UniAM, an innovative Compressive Attention Matching framework. What distinguishes UniAM is its unique capability to exploit the self-attention mechanism inherent in ViT, allowing it to adeptly capture the most pertinent information necessary for Universal Domain Adaptation. This is further complemented by its innovative compressive reconstruction module and residual-based transferability criterion, which together enable effective domain alignment. It's worth noting that UniAM stands as a pioneering method that directly harnesses the attention capabilities of vision transformers, specifically for classification tasks. Through extensive experiments on four benchmark datasets, we've found that our approach consistently eclipses the state-of-the-art UniDA method in both common set accuracy and "unknown" class accuracy. We hope these findings will provide a new perspective for domain adaptation and other fields such as out of distribution detection in this future.

# References

[1] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.

[2] Silvia Bucci, Mohammad Reza Loghmani, and Tatiana Tommasi. On the effectiveness of image rotation for open set domain adaptation. In *European Conference on Computer Vision*, pages 422–438. Springer, 2020.

[3] Zhangjie Cao, Mingsheng Long, Jianmin Wang, and Michael I Jordan. Partial transfer learning with selective adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2724–2732, 2018.

[4] Zhangjie Cao, Lijia Ma, Mingsheng Long, and Jianmin Wang. Partial adversarial domain adaptation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 135–150, 2018.

[5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.

[6] Wanxing Chang, Ye Shi, Hoang Duong Tuan, and Jingya Wang. Unified optimal transport framework for universal domain adaptation. *arXiv preprint arXiv:2210.17067*, 2022.

[7] Liang Chen, Qianjin Du, Yihang Lou, Jianzhong He, Tao Bai, and Minghua Deng. Mutual nearest neighbor contrast and hybrid prototype self-training for universal domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 6248–6257, 2022.

[8] Liang Chen, Yihang Lou, Jianzhong He, Tao Bai, and Minghua Deng. Evidential neighborhood contrastive learning for universal domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 6258–6267, 2022.

[9] Liang Chen, Yihang Lou, Jianzhong He, Tao Bai, and Minghua Deng. Geometric anchor correspondence mining with uncertainty modeling for universal domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16134–16143, 2022.

[10] Yi-Chen Chen, Vishal M Patel, P Jonathon Phillips, and Rama Chellappa. Dictionary-based face recognition from video. In *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part VI 12*, pages 766–779. Springer, 2012.

[11] Yuejie Chi and Fatih Porikli. Classification and boosting with multiple collaborative representations. *IEEE transactions on pattern analysis and machine intelligence*, 36(8):1519–1531, 2013.

[12] Jia Deng. A large-scale hierarchical image database. *Proc. of IEEE Computer Vision and Pattern Recognition, 2009*, 2009.

[13] Bo Fu, Zhangjie Cao, Mingsheng Long, and Jianmin Wang. Learning to detect open classes for universal domain adaptation. In *European Conference on Computer Vision*, pages 567–583. Springer, 2020.

[14] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015.

[15] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016.

[16] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018.

[17] Rohit Girdhar, Joao Carreira, Carl Doersch, and Andrew Zisserman. Video action transformer network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 244–253, 2019.

[18] Haitao He, Haoran Niu, Jianzhou Feng, Qian Wang, and Qikai Wei. A prototype network enhanced relation semantic representation for few-shot relation extraction. *Human-Centric Intelligent Systems*, 3(1):1–12, 2023.

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[20] Katherine Hermann, Ting Chen, and Simon Kornblith. The origins and prevalence of texture bias in convolutional neural networks. *Advances in Neural Information Processing Systems*, 33:19000–19015, 2020.

[21] Ji Huang, Minbo Ma, Yongsheng Dai, Jie Hu, and Shengdong Du. Dbaformer: A double-branch attention transformer for long-term time series forecasting. *Human-Centric Intelligent Systems*, pages 1–12, 2023.

[22] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.

[23] Vasima Khan and Tariq Azfar Meenai. Pretrained natural language processing model for intent recognition (bert-ir). *Human-Centric Intelligent Systems*, 1(3-4):66–74, 2021.

[24] Alexander Kolesnikov, Alexey Dosovitskiy, Dirk Weissenborn, Georg Heigold, Jakob Uszkoreit, Lucas Beyer, Matthias Minderer, Mostafa Dehghani, Neil Houlsby, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. 2021.

[25] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.

[26] Jogendra Nath Kundu, Suvaansh Bhambri, Akshay Kulkarni, Hiran Sarkar, Varun Jampani, and R Venkatesh Babu. Subsidiary prototype alignment for universal domain adaptation. *arXiv preprint arXiv:2210.15909*, 2022.

[27] Jogendra Nath Kundu, Naveen Venkat, R Venkatesh Babu, et al. Universal source-free domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4544–4553, 2020.

[28] Barbara Landau, Linda B Smith, and Susan S Jones. The importance of shape in early lexical learning. *Cognitive development*, 3(3):299–321, 1988.

[29] Guangrui Li, Guoliang Kang, Yi Zhu, Yunchao Wei, and Yi Yang. Domain consensus clustering for universal domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9757–9766, 2021.

[30] Jian Liang, Yunbo Wang, Dapeng Hu, Ran He, and Jiashi Feng. A balanced and uncertainty-aware approach for partial domain adaptation. In *European Conference on Computer Vision*, pages 123–140. Springer, 2020.

[31] Omri Lifshitz and Lior Wolf. A sample selection approach for universal domain adaptation. *arXiv preprint arXiv:2001.05071*, 2020.

[32] Hong Liu, Zhangjie Cao, Mingsheng Long, Jianmin Wang, and Qiang Yang. Separate to adapt: Open set domain adaptation via progressive separation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2927–2936, 2019.

[33] Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, and Zheyan Shen. Heterogeneous risk minimization. In *International Conference on Machine Learning*, pages 6804–6814. PMLR, 2021.

[34] Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, and Zheyan Shen. Kernelized heterogeneous risk minimization. *arXiv preprint arXiv:2110.12425*, 2021.

[35] Xiaofeng Liu, Zhaofeng Li, Lingsheng Kong, Zhihui Diao, Junliang Yan, Yang Zou, Chao Yang, Ping Jia, and Jane You. A joint optimization framework of low-dimensional projection and collaborative representation for discriminative classification. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pages 1493–1498. IEEE, 2018.

[36] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.

[37] Shuang Luo, Yinchuan Li, Jiahui Li, Kun Kuang, Furui Liu, Yunfeng Shao, and Chao Wu. S2rl: Do we really need to perceive all states in deep multi-agent reinforcement learning? In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1183–1191, 2022.

[38] Zheqi Lv, Zhengyu Chen, Shengyu Zhang, Kun Kuang, Wenqiao Zhang, Mengze Li, Beng Chin Ooi, and Fei Wu. Ideal: Toward high-efficiency device-cloud collaborative and dynamic recommendation system. *arXiv preprint arXiv:2302.07335*, 2023.

[39] Zheqi Lv, Wenqiao Zhang, Shengyu Zhang, Kun Kuang, Feng Wang, Yongwei Wang, Zhengyu Chen, Tao Shen, Hongxia Yang, Beng Chin Ooi, et al. Duet: A tuning-free device-cloud collaborative parameters generation framework for efficient device model generalization. In *Proceedings of the ACM Web Conference 2023*, pages 3077–3085, 2023.

[40] Xu Ma, Junkun Yuan, Yen-wei Chen, Ruofeng Tong, and Lanfen Lin. Attention-based cross-layer domain alignment for unsupervised domain adaptation. *Neurocomputing*, 499:1–10, 2022.

[41] Muhammad Muzammal Naseer, Kanchana Ranasinghe, Salman H Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers. *Advances in Neural Information Processing Systems*, 34:23296–23308, 2021.

[42] Daniel Neimark, Omri Bar, Maya Zohar, and Dotan Asselmann. Video transformer network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3163–3172, 2021.

[43] Pau Panareda Busto and Juergen Gall. Open set domain adaptation. In *Proceedings of the IEEE international conference on computer vision*, pages 754–763, 2017.

[44] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1406–1415, 2019.

[45] Xingchao Peng, Ben Usman, Neela Kaushik, Dequan Wang, Judy Hoffman, and Kate Saenko. Visda: A synthetic-to-real benchmark for visual domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 2021–2026, 2018.

[46] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *European conference on computer vision*, pages 213–226. Springer, 2010.

[47] Kuniaki Saito, Donghyun Kim, Stan Sclaroff, and Kate Saenko. Universal domain adaptation through self supervision. *Advances in neural information processing systems*, 33:16282–16292, 2020.

[48] Kuniaki Saito and Kate Saenko. Ovanet: One-vs-all network for universal domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9000–9009, 2021.

[49] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3723–3732, 2018.

[50] Kuniaki Saito, Shohei Yamamoto, Yoshitaka Ushiku, and Tatsuya Harada. Open set domain adaptation by backpropagation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 153–168, 2018.

[51] Tao Shen, Jie Zhang, Xinkang Jia, Fengda Zhang, Gang Huang, Pan Zhou, Kun Kuang, Fei Wu, and Chao Wu. Federated mutual learning. *arXiv preprint arXiv:2006.16765*, 2020.

[52] Zheyan Shen, Jiashuo Liu, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*, 2021.

[53] Tao Sun, Cheng Lu, Tianshuo Zhang, and Haibin Ling. Safe self-refinement for transformer-based domain adaptation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7191–7200, 2022.

[54] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.

[55] Shikhar Tuli, Ishita Dasgupta, Erin Grant, and Thomas L Griffiths. Are convolutional neural networks or transformers more like human vision? *arXiv preprint arXiv:2105.07197*, 2021.

[56] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017.

[57] Vladimir Vapnik. Principles of risk minimization for learning theory. *Advances in neural information processing systems*, 4, 1991.

[58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[59] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5018–5027, 2017.

[60] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021.

[61] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8741–8750, 2021.

[62] John Wright, Allen Y Yang, Arvind Ganesh, S Shankar Sastry, and Yi Ma. Robust face recognition via sparse representation. *IEEE transactions on pattern analysis and machine intelligence*, 31(2):210–227, 2008.

[63] Yang Wu, Vansteenberge Jarich, Masayuki Mukunoki, and Michihiko Minoh. Collaborative representation for classification, sparse or non-sparse? *arXiv preprint arXiv:1403.1353*, 2014.

[64] Tongkun Xu, Weihua Chen, Pichao Wang, Fan Wang, Hao Li, and Rong Jin. Cdtrans: Cross-domain transformer for unsupervised domain adaptation. *arXiv preprint arXiv:2109.06165*, 2021.

[65] Jinyu Yang, Jingjing Liu, Ning Xu, and Junzhou Huang. Tvt: Transferable vision transformer for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 520–530, 2023.

[66] Kaichao You, Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Universal domain adaptation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2720–2729, 2019.

[67] Junkun Yuan, Xu Ma, Defang Chen, Kun Kuang, Fei Wu, and Lanfen Lin. Domain-specific bias filtering for single labeled domain generalization. *International Journal of Computer Vision*, 131(2):552–571, 2023.

[68] Junkun Yuan, Xu Ma, Ruoxuan Xiong, Mingming Gong, Xiangyu Liu, Fei Wu, Lanfen Lin, and Kun Kuang. Instrumental variable-driven domain generalization with unobserved confounders. *ACM Transactions on Knowledge Discovery from Data*, 2023.

[69] Fengda Zhang, Kun Kuang, Long Chen, Yuxuan Liu, Chao Wu, and Jun Xiao. Fairness-aware contrastive learning with partially annotated sensitive attributes. In *The Eleventh International Conference on Learning Representations*, 2022.

[70] Fengda Zhang, Kun Kuang, Zhaoyang You, Tao Shen, Jun Xiao, Yin Zhang, Chao Wu, Yueting Zhuang, and Xiaolin Li. Federated unsupervised representation learning. *arXiv preprint arXiv:2010.08982*, 2020.

[71] Lei Zhang, Meng Yang, and Xiangchu Feng. Sparse representation or collaborative representation: Which helps face recognition? In *2011 International conference on computer vision*, pages 471–478. IEEE, 2011.

[72] Min Zhang, Siteng Huang, and Donglin Wang. Domain generalized few-shot image classification via meta regularization network. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3748–3752. IEEE, 2022.

[73] Min Zhang, Donglin Wang, and Sibo Gai. Knowledge distillation for model-agnostic meta-learning. In *ECAI 2020*, pages 1355–1362. IOS Press, 2020.

[74] Min Zhang, Zifeng Zhuang, Zhitao Wang, Donglin Wang, and Wenbin Li. Rotogbml: Towards out-of-distribution generalization for gradient-based meta-learning. *arXiv preprint arXiv:2303.06679*, 2023.

[75] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021.

[76] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020.