# MAPS: A Noise-Robust Progressive Learning Approach for Source-Free Domain Adaptive Keypoint Detection

Yuhe Ding, Jian Liang, Bo Jiang, Aihua Zheng and Ran He, *Senior member, IEEE*

arXiv:2302.04589v1 [cs.CV] 9 Feb 2023

*Abstract*—Existing cross-domain keypoint detection methods always require accessing the source data during adaptation, which may violate the data privacy law and pose serious security concerns. Instead, this paper considers a realistic problem setting called source-free domain adaptive keypoint detection, where only the well-trained source model is provided to the target domain. For the challenging problem, we first construct a teacher-student learning baseline by stabilizing the predictions under data augmentation and network ensembles. Built on this, we further propose a unified approach, Mixup Augmentation and Progressive Selection (MAPS), to fully exploit the noisy pseudo labels of unlabeled target data during training. On the one hand, MAPS regularizes the model to favor simple linear behavior in-between the target samples via self-mixup augmentation, preventing the model from over-fitting to noisy predictions. On the other hand, MAPS employs the self-paced learning paradigm and progressively selects pseudo-labeled samples from 'easy' to 'hard' into the training process to reduce noise accumulation. Results on four keypoint detection datasets show that MAPS outperforms the baseline and achieves comparable or even better results in comparison to previous non-source-free counterparts. The code will be released at **https://github.com/YuheD/MAPS**.

*Index Terms*—Source-free Domain Adaptation, Keypoint Detection, Noise-Robust Learning.

## I. INTRODUCTION

Unsupervised domain adaptation (UDA), a paradigm that aims to transfer knowledge from a label-rich source domain to another unlabeled target domain, is proposed to reduce the burden of manual annotations [1]. Typically, existing UDA methods resort to feature alignment [2], [3] or domain translation [4] so that the discriminative model learned on the labeled source domain could perform well in the unlabeled target domain. However, these methods always require access to the source data when adapting to the target domain, which may be infeasible due to the growing restrictions on data privacy and storage. To tackle this issue, a line of recent works [5]–[8] consider a source data-free domain adaptation (SFDA) scheme in which only the well-trained classification model from the source domain is provided for the target domain. Besides the widely-studied image classification task [5], the

Yuhe Ding and Bo Jiang are with the School of Computer Science and Technology, Anhui University. E-mail: madao3c@foxmail.com; jiangbo@ahu.edu.cn.

Aihua Zheng is with the School of Artificial Intelligence, Anhui University. E-mail: ahzheng214@foxmail.com.

Jian Liang and Ran He are with the National Laboratory of Pattern Recognition (NLPR) and the Center for Research on Intelligent Perception and Computing, CASIA. E-mail: liangjian92@gmail.com; rhe@nlpr.ia.ac.cn.

Jian Liang is the corresponding author.

SFDA scheme has also been extended to other visual tasks, *e.g.*, semantic segmentation [8]–[11], object detection [12]–[14] and point cloud recognition [7], [15]. Though much effort has been devoted to the classification problem, only a few works [16] explore the SFDA scheme for regression problems.

This paper for the first time studies the SFDA scheme in a popular regression task—image keypoint detection—which aims to localize the predefined object keypoints in 2D images [17]. Prior domain adaptive keypoint detection works typically align the source and target domains by elaborately designed alignment strategies, *e.g.*, domain adversarial training [18], [19] and style transfer model [20]. However, all of them require the existence of source data and target data at the same time. Besides, current SFDA methods for classification always enlarge the margins of boundaries between different classes on the target domain, *e.g.*, pseudo-labeling [5], [7], information maximization [5], contrastive learning [6]. Such techniques are also not applicable to the keypoint detection task since the regression space is continuous without a clear decision boundary.

For the challenging SFDA keypoint detection task, we come up with a simple baseline method by exploiting the mean-teacher framework [21], where the teacher model is an exponential moving average (EMA) of the student network. To ensure consistency across models, target samples with different augmentations are utilized in the teacher and student models, and the mean squared error between their predictions is minimized. This method only fine-tunes the model with the consistency loss and no longer requires access to the source data. Nevertheless, it does not fully explore the semantic information in the unlabeled target domain, leading to sub-optimal adaptation performance. Therefore, we resort to pseudo-labeling [22], a popular semi-supervised learning technique that generates the pseudo labels for unlabeled data. Since there is no labeled data, solely relying on these noisy pseudo labels is inevitably harmful due to noise accumulation, misleading the model over time.

In this paper, we propose a new approach termed Mixup Augmentation and Progressive Selection (MAPS) by alleviating the noise accumulation issue for source-free keypoint detection. Firstly, MAPS mixes the target samples with their shuffle version to construct the self-mixup augmentation [23]. This manner regularizes the model to produce outputs that vary linearly with the inputs, preventing over-fitting to outliers and thereby improving the model's robustness. Secondly, to explore reliable pseudo labels, MAPS adopts the self-paced

learning scheme (SPL) [24] to alternately select easy samples with high confidence and apply a standard regression loss on them. In fact, the number of selected samples is governed by a dynamic weight that is annealed until the training loss converges. Prioritizing learning about easy samples is proven to prevent the model from getting stuck in a bad local optimum [25]. These two strategies respectively consider the convex behavior among samples and the quality of the pseudo label per sample, making them complement each other. To validate the effectiveness of our method, we conduct three experiments on human, hand, and animal keypoint detection. The main contributions of our work are summarized as follows:

- We formulate a realistic and challenging task, *source-free domain adaptive keypoint detection*, which is among the first to study cross-domain regression problems without the source data.
- Built on the devised baseline via mean-teacher, we propose a new approach termed Mixup Augmentation and Progressive Selection (MAPS) which fully exploits the noisy pseudo labels in the target domain.
- Experiments on various keypoint detection benchmarks demonstrate that MAPS outperforms the baseline and achieves results comparable, or even better than, existing non-source-free counterparts.

## II. RELATED WORK

### A. Source-free Domain Adaptation

As a typical problem in transfer learning [1], [26], [27], unsupervised domain adaptation (UDA) [1]–[3], [28] aims to improve the performance of the model on an unlabeled target domain with the aid of label-rich source domain. However, as data privacy has received increasing attention in recent years, source-free domain adaptation (SFDA) [5] is proposed to solve the UDA problem with only a well-trained source model.

There are two primary types of SFDA approaches, *i.e.*, generation-based, and self-training-based. Generation-based methods [29]–[33] attempt to restore the source features or images to compensate for the absence of source data. Self-training-based methods [5], [6], [15], [34], [35] leverage self-supervised techniques to explore the intrinsic structure of the target domain. In recent years, there are also some practices [32], [36] to solve this problem by selecting part of the target data as a pseudo source domain, to compensate for the unseen source domain. The following works develop kinds of variants of SFDA, *e.g.*, active SFDA [37], black-box SFDA [38], multi-source SFDA [39], semi-supervised SFDA [40]. In reality, the vanilla SFDA along with these variants is mainly applied in multiple visual classification tasks, *e.g.*, object recognition [5], semantic segmentation [8], [10], and object detection [12]. The SFDA scheme for classification always enlarges the margins of decision boundaries between different classes on the target domain, thereby improving the generalization performance. However, the regression space is continuous without a clear boundary, leading to these methods cannot apply to regression tasks directly. There are a few works proposed for measurement shift [41], blind image quality assessment [16]. For the general keypoint detection task, it

is still blank. To the best of our knowledge, we are the first to investigate the SFDA scheme on the keypoint detection task.

### B. Domain Adaptive Keypoint Detection

Keypoint detection is a fundamental visual problem. It is laborious and time-consuming to collect the labeled data for traditional keypoint detection methods [42]–[46]. Domain adaptive keypoint detection [19], [47]–[52] attempts to address this issue by transferring information from a labeled source domain to an unlabeled target domain, avoiding the extra annotation costs. These methods are primarily divided into 3D and 2D keypoint detection. In 3D keypoint detection, this problem is complex and there are many prior works. For instance, [51] proves that the neural networks can perform well when the data is pre-processed to extract cues about the person's motion, notably as optical flow and the motion of 2D keypoints, and therefore propose to use motion as a simple way to bridge a Synthetic-to-realistic gap when the video is available. UDA-COPE [52] introduces a bidirectional filtering method between the predicted normalized object coordinate space (NOCS) map and observed point cloud, to promote teacher-student consistency training.

This paper focuses on another type, *i.e.*, 2D keypoint detection. In this field, existing methods can be primarily divided into human keypoint detection and animal keypoint detection. For human keypoint detection, RegDA [19], and MarsDA [53] introduce adversarial regressors to narrow the domain gap. C-GAC [47] integrates the proposed prediction confidence into self-training to obtain reliable pseudo labels. For animal keypoint detection, CC-SSL [54], and UDA-Animal [18] are based on transformation consistency and the pseudo-label refinery technique. Besides, unified domain adaptive pose estimation (UDAPE) [20] proposes to align representations using both input-level and output-level cues, and provides a unified framework for both human keypoint detection and animal keypoint detection problems. Generally, the source data is essential for these carefully designed alignment strategies, which may violate data privacy law and pose serious concerns, while our method solves this problem under a source data-free setting.

### C. Curriculum Learning

Curriculum learning (CL) [55] is a training strategy that trains a model from easy to hard, which simulates the learning principle of humans and animals. CL serves as a general training strategy and has been used in multiple applications of CL in machine learning [56]–[59]. In these applications, the motivations can be categorized for applying CL into two groups: to guide, regularizing the training towards better regions in parameter space (with steeper gradients) from the perspective of the optimization problem, and to denoise, focusing on the high-confidence easier area to alleviate the interference of noisy data as from the perspective of data distribution. A general CL framework has a Difficulty Measurer, which decides the relative "easiness" of each data example, and a Training Scheduler, which decides the sequence of data subsets throughout the training process based on the judgment
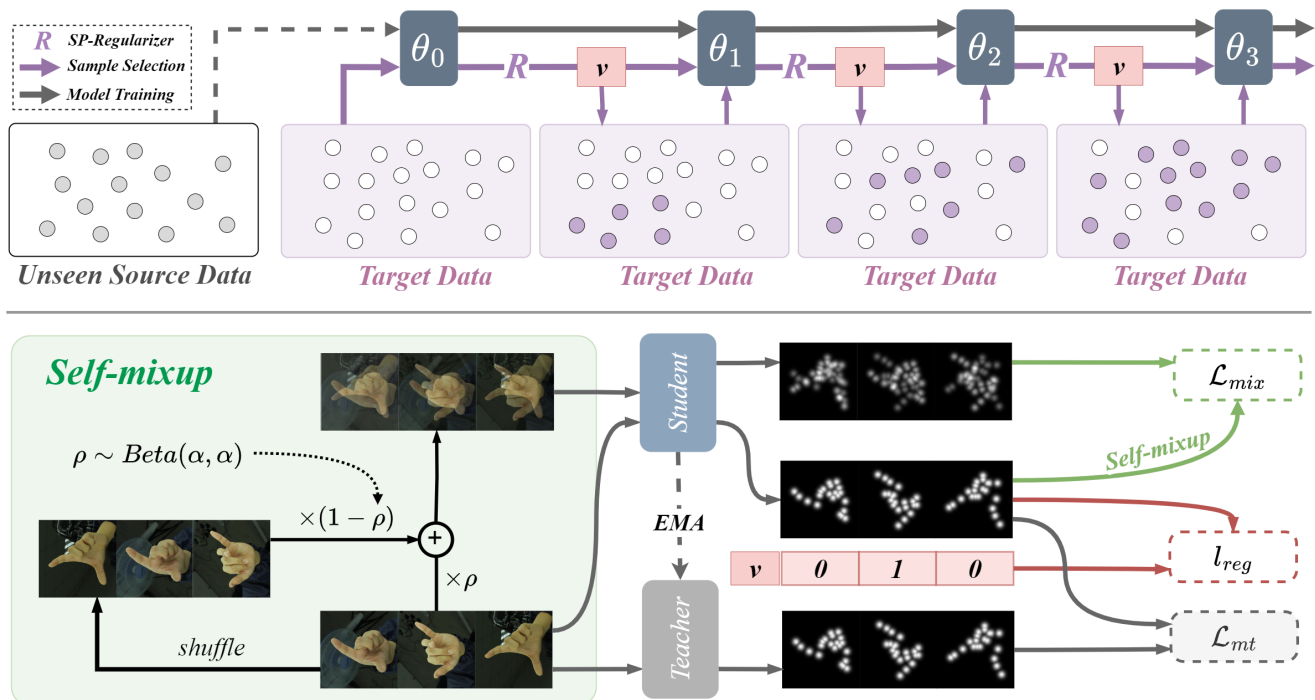
Fig. 1: The pipeline of Mixup Augmentation and Progressively Selection (MAPS). Given the well-trained source model $\theta_0$, MAPS progressively selects easy samples from target samples, *i.e.*, assigns weight, 1 for easy samples and 0 for the remaining samples, the weights are notated as symbol $v$. Then these samples participate in the model updating with self-mixup training, which regularizes the model to favor a linear behavior in-between the target samples. These two aspects are alternatively performed until the loss function converges.

from the Difficulty Measurer. Based on this framework, CL is primarily divided into two cases, *i.e.*, predefined CL and automatic CL. The predefined CL case measures the difficulty of samples and schedules the training process based on prior knowledge, and the automatic CL case adopts a data-driven method to learn any (or both) aspects. In predefined CL, researchers manually design various Difficulty Measurers mainly based on the data characteristics of specific tasks. For instance, in NLP tasks, sentence length is the most popular Difficulty Measure in NLP tasks [60]. This is because researchers intuitively think the complexity of a sentence or paragraph can be expressed by sentence length.

Self-paced learning (SPL) [24], [61] is a typical example of automatic curriculum learning. The easiest samples with the lowest training loss are used to train the models, this portion of the easiest samples gradually increases according to a specified scheduler. Based on the scheduler, the models are able to learn at their own pace, including deciding what, how, when, and how long to study. Prior work formulates the critical principle of SPL as a concise model [24], and proves that the SPL regime is convergent by adopting an alternative optimization strategy (AOS) [62]. Generally, previous self-paced learning methods focus on the supervised learning scenario, and no few works study the application of SPL on unsupervised learning tasks. Traditional supervised SPL achieves promising results, while unsupervised SPL still remains much room for improvement. In this paper, we study unsupervised SPL in regression problems.

## III. METHOD

We define the notations in Section III-A and generate the source model by a standard regression loss in section III-B. Then we construct a simple baseline by the mean-teacher framework in Section III-C and introduce the proposed noise-robustness progressive learning method in Section III-D.

### A. Notations

For a standard unsupervised domain adaptive 2D keypoint detection, we have $n_s$ labeled samples $\{x_i^s, y_i^s\}_{i=1}^{n_s}$ from the source domain $D_s = \{\mathcal{X}_s, \mathcal{Y}_s^K\}$, and $n_t$ unlabeled data $\{x_i^t\}_{i=1}^{n_t}$ from the target domain $D_t = \{\mathcal{X}_t\}$, where $\mathcal{X}_s, \mathcal{X}_t \in \mathcal{R}^{H \times W \times 3}$ are the input space, $\mathcal{Y}_s^K \in \mathcal{R}^2$ is the output space, $K$ is the number of keypoints in each input image. The goal is to find a keypoint detector to predict the labels $\{y_i^t\}_{i=1}^{n_t}$ of the target domain.

The source-free unsupervised domain adaptive keypoint detection task aims to learn a target keypoint detector $f_t : \mathcal{X}_t \to \mathcal{Y}_t$ and infer $\{y_i^t\}_{i=1}^{n_t}$ with only $\{x_i^t\}_{i=1}^{n_t}$ and the source keypoint detector $f_s$. We first generate the source keypoint detection model by a standard regression method, and then we transfer the model to the target domain without accessing the source data.

### B. Source Model Generation

We consider to train the source keypoint detection model $f_s : \mathcal{X}_s \to \mathcal{Y}_s$ by minimizing the following standard super-

vised regression loss,

$$\mathcal{L}_{src}(\mathcal{X}_s, \mathcal{Y}_s) = \sum_{i=1}^{n_s} \|f_s(x_i^s) - H(y_i^s)\|_2, \qquad (1)$$

where $H(\cdot)$ denotes the heatmap generation function [63]. To further improve the performance of the source model and facilitate the following target data alignment, we use the random data augmentation technique [18] at this stage. Taking the augmentation into consideration, the objective function is reformulated as:

$$\mathcal{L}_{src}^{aug}(\mathcal{X}_s, \mathcal{Y}_s) = \sum_{i=1}^{n_s} \|f_s(\mathcal{A}^s(x_i^s)) - \mathcal{A}^s(H(y_i^s))\|_2, \qquad (2)$$

where $\mathcal{A}^s$ denotes the augmentation function.

### C. Teacher-Student Consistency Learning

Given the well-trained source model $f_s$, we first construct a baseline method based on the mean-teacher framework [18], [20], [21]. Specifically, the student model $f_t$ and teacher model $f_t'$ have an identical architecture with the given source model $f_s$, and the teacher parameter $\theta'$ is updated with the exponential moving average (EMA) of the student parameter $\theta$:

$$\theta_t' = \eta\theta_{t-1}' + (1 - \eta)\theta_t, \qquad (3)$$

where $t$ denotes the step of training, and $\eta$ denotes the smoothing coefficient which is set to 0.999 by default.

Here we also introduce the random data augmentation function $A^t$ on target data, the input of student $\mathcal{A}_1^t(x_i^t)$ and the input of teacher $\mathcal{A}_2^t(x_i^t)$ are generated from different augmentation $\mathcal{A}_1^t$ and $\mathcal{A}_2^t$, sampled from $\mathcal{A}^t$. Besides, considering the missing and occluded keypoints in some images, we only force the teacher-student consistency on the predicted points with the maximum activation greater than a threshold $\tau$. The object function is formulated as:

$$\mathcal{L}_{mt}(\mathcal{X}_t) = \sum_{i=1}^{n_t} \sum_{k=1}^{K} \mathbb{1}(h_k' \geq \tau)\|\tilde{\mathcal{A}}_2^t(h_k') - \tilde{\mathcal{A}}_1^t(h_k))\|_2, \quad (4)$$

where $\tilde{A}_1^t$ and $\tilde{A}_2^t$ denote the inverse function of $\mathcal{A}_1^t$ and $\mathcal{A}_2^t$, $\{h_k\}_{k=1}^K = f_t(\mathcal{A}_1^t(x_i^t))$ denote the outputs of student model, $\{h_k'\}_{k=1}^K = H(f_t'(\mathcal{A}_2^t(x_i^t))$ denote the normalized outputs of teacher model, $\mathbb{1}(\cdot)$ denotes the indicator function, and $K$ denotes the number of keypoints.

### D. Noise-Robust Progressive Learning

Based on the baseline depicted in section III-C, we further employ a simple strategy, pseudo-labeling [22], which directly uses the current prediction of the model as the ground truth, to explore the semantic information in the target domain. However, this manner inevitably introduces many erroneous labels, making the model suffer from noise accumulation, which obstructs true signal and biasing estimation of corresponding parameters [64], leading to sub-optimal results. As shown in Fig. 1, we further propose Mixup Augmentation and Progressive Selection (MAPS), to avoid the over-confidence

in outliers and select reliable pseudo labels progressively to participate in the model training.

**Mixup Augmentation.** Based on the baseline model, we propose self-mixup augmentation to favor simple linear behavior in-between target data during training, thereby preventing confirmation bias [23] and improving the model's robustness. To be specific, we mix the augmented target samples $\mathcal{A}^t(x_i^t)$ and another random sample $\mathcal{A}^t(x_j^t)$ $(i \neq j)$ in the current mini-batch, with a mix ratio $\rho \sim Beta(\alpha, \alpha)$ to construct the self-mixup augmentation, and then we introduce the self-mixup loss in the following:

$$\begin{aligned} \mathcal{L}_{mix}(\mathcal{X}_t) &= \sum_{i=1}^{n_t} \|f_t(x_i^m) - h_i^m\|_2, \\ x_i^m &= \rho\mathcal{A}^t(x_i^t) + (1 - \rho)\mathcal{A}^t(x_j^t), \\ h_i^m &= \rho f_t(\mathcal{A}^t(x_i^t)) + (1 - \rho)f_t(\mathcal{A}^t(x_j^t)). \end{aligned} \qquad (5)$$

The self-mixup loss regularizes the model to favor simple linear behavior in-between the target samples, preventing the over-fitting of outliers.

**Progressive Selection.** The pseudo-labeling technique [22] picks up the current prediction as used as if they were true labels. However, it is inevitable to introduce noisy labels, which might mislead the model over time. We introduce the self-paced learning regime [24] to progressively select samples with high-confident pseudo labels, and the pseudo-label regression loss is imposed on them. This aim can be realized by minimizing the following objective:

$$\mathcal{L}_{spl}(\mathcal{X}_t) = \sum_{i=1}^{n_t} v_i l_{reg}(x_i^t) + R(v_i, \lambda), \ s.t. \ v_i \in [0, 1] \quad (6)$$

where $\lambda$ denotes the age parameter for adjusting the learning pace, $v_i$ denotes the weight for the $i$-th sample, $R(v_i, \lambda)$ denotes the self-paced regularizer (SP-regularizer), and $l_{spl}$ denotes the self-paced loss weighted by $v_i$. The self-paced loss for each sample $x_i^t$ is defined by the mean squared error, *i.e.*, pseudo-label regression loss:

$$l_{reg}(x_i^t) = \|f_t(\mathcal{A}^t(x_i^t)) - \mathcal{A}^t(H(\hat{y}_i^t))\|_2, \qquad (7)$$

where $\hat{y}_i^t$ is the pseudo label of $x_i^t$, which is the 2D coordinate of the maximum activation of the student output $f_t(x_i^t)$.

Here we adopt the hard SP-regularizer [24] to solve the sample weights. The formulation and the closed-form solution are below:

$$R(v_i, \lambda) = -\lambda v_i; \ v_i = \begin{cases} 1, & l_{reg}(x_i^t) < \lambda \\ 0, & l_{reg}(x_i^t) \geq \lambda \end{cases}. \qquad (8)$$

The model parameters are fixed during selection, therefore, $l_{reg}$ can be seen as a difficulty measurer, and the age parameter $\lambda$ is a threshold. This regime progressively selects the easy samples by assigning the weight 1 for them.

The age parameter $\lambda$ is updated by a baby step scheduler $S$ [65]. To be specific, we pre-defines an increasing sequence $M = \{M_1, M_2, M_3, ..., M_Q\}$. In the $q$-th round, we select $M_q$ samples, and $\lambda$ can be calculated by:

$$\lambda = S\big(M_q, \{l_{reg}(x_i^t)\}_{i=1}^{n_t}\big), \qquad (9)$$

where $S(m, L)$ denotes the $m$-th smallest value of loss function $L$.

**Overall Objective.** With the above two aspects, the overall objective of MAPS can be summarized as:

$$\min_{\theta, v} \mathcal{L}_{mt}(\mathcal{X}_t) + \beta_m \mathcal{L}_{mix}(\mathcal{X}_t) + \beta_s \mathcal{L}_{spl}(\mathcal{X}_t), \quad (10)$$

where $\beta_s$, $\beta_m$ are hyper-parameters to control the weights of the corresponding loss terms. We adopt the alternative optimization strategy to optimize the objective function in the following steps.

1) Select easy samples, *i.e.*, solve weights $v$, when $\theta$ and $\lambda$ are fixed. Optimizing Eq. (10) is equal to solving the following problem:

$$\min_{v} \mathcal{L}_{spl}(\mathcal{X}_t), \ s.t. \ v_i \in [0, 1], \quad (11)$$

we adopt the hard SP-regularizer here, and the closed-form solution of $v$ is shown in Eq. (8).

2) Train the model, *i.e.*, update $\theta$, when $v$ and $\lambda$ are fixed. Optimizing Eq. (10) is equal to solving the following problem:

$$\min_{\theta} \mathcal{L}_{mt}(\mathcal{X}_t) + \beta_m \mathcal{L}_{mix}(\mathcal{X}_t) + \beta_s \sum_{i=1}^{n_t} v_i l_{reg}(x_i^t). \quad (12)$$

In addition, according to the given baby step scheduler in Eq. (9), we also update the age parameter $\lambda$ to include more training samples during the optimization process. The overall algorithm is shown in Algorithm 1.

---

**Algorithm 1** Algorithm of the proposed MAPS.

**Input:** Target dataset $\mathcal{X}_t$; a well-trained source model $f_s$;
**Output:** New target model $f_t$ and $f_t'$;
 1: Initialize the student $f_t$ and teacher $f_t'$ with $f_s$;
 2: **repeat**
 3:   Fix $f_t$, update $v$ by Eq. (11); // sample selection
 4:   Update $f_t$ by Eq. (12); // model updating
 5:   Update $f_t'$ by Eq. (3);
 6:   Update $\lambda$ based on $S$ by Eq. (9);
 7: **until** Progressive selection rounds are exhausted.

---

## IV. EXPERIMENT

We first give the experimental setup in Section IV-A, then we conduct the experiments and analyze the results on hand keypoint detection, human keypoint detection, and animal keypoint detection in Section IV-B, IV-C, IV-D, respectively. To validate the effectiveness of the proposed modules, we conduct the ablation study in Section IV-E and analyze the sensitivity in Section IV-F.

### A. Setup

The architecture of our model is based on Simple Baseline [44] with the backbone of pre-trained ResNet101 [66]. Following prior works [18], [20], [54], we add texture and geometric augmentation including Gaussian noise, Gaussian blurring, rotation, and random 2D translation for both source model generation and target model training. The parameter $\alpha$ of the beta distribution is set to 0.75, and the weights of

TABLE I: PCK@0.05 on task **RHD→H3D**. (Best value is in Redorange color. Second best value is in periwinkle color.)

| Method | SF | MCP | PIP | DIP | Fin | Avg. |
|---|---|---|---|---|---|---|
| CCSSL [54] | × | 81.5 | 79.9 | 74.4 | 64.0 | 75.1 |
| UDA-Animal [18] | × | 82.3 | 79.6 | 72.3 | 61.5 | 74.1 |
| RegDA [19] | × | 79.6 | 74.4 | 71.2 | 62.9 | 72.5 |
| UDAPE [20] | × | 86.7 | 84.6 | 78.9 | 68.1 | 79.6 |
| Source only | - | 65.0 | 62.4 | 61.1 | 54.0 | 60.1 |
| MT | ✓ | 85.0 | 84.6 | 77.8 | 67.3 | 78.8 |
| MAPS | ✓ | 86.9 | 84.8 | 79.1 | 69.3 | 80.0 |

loss functions $\beta_m$ and $\beta_s$ are both empirically set to 1. In the source model generation stage, we adopt Adam [67] as the optimizer and set the learning rate as 2e-4, and it decreased to 2e-5 at 15,000 steps and 2e-6 at 20,000 steps, while there is a total of 25,000 steps in this stage. In the target model training, we also adopt the Adam optimizer and learning rate 2e-4, and the learning rate decreases to 1e-4 at 2,500 steps with a total of 7,500 steps. Our basic augmentation is based on UDAPE [20], which adds rotation and random 2D translation on the augmentation in RegDA [19]. We use the Percentage of Correct Keypoints (PCK) as our metric, in which estimation is considered correct if its distance from the ground truth is less than a fraction of 0.05 of the image size. We randomly run our methods three times with different random seeds $\{0,1,2\}$ via **PyTorch** and report the average accuracies.

We conduct experiments on three tasks with four different target datasets, and compare MAPS with several state-of-the-art *non-source-free* methods, including the animal keypoint detection methods CCSSL [54] and UDA-Animal [18]; the human and hand keypoint detection method RegDA [19], and the unified method UDAPE [20]. In the table of comparative experiments, 'source only' denotes using the source model for prediction, 'MT' is our baseline method introduced in section III-C, and MAPS is the full version of our method. All the results except ours are from the original paper of UDAPE, and the dataset split policy follows UDAPE and RegDA.

### B. Hand Keypoint Detection

**Dataset.** Rendered Hand Pose Dataset [68] (**RHD**) is a synthetic dataset containing 41,258 training images and 2,728 testing images with corresponding 21 hand keypoints labels. Hand-3D-Studio [69] (**H3D**) is a real-world multi-view hand image dataset containing 22k images in total. Following the split way in RegDA [19], we select 3.2k images as the testing set, and the remaining as the training set.

**Implementation Details.** We use **RHD** as the source domain and **H3D** as the target domain. The number of progressive selection rounds $Q$ and the increasing sequence $M$ in the baby step schedular are set to 3 and $\{0.25N, 0.35N, 0.45N\}$, respectively, where $N$ is the total number of training samples. We report 21 keypoints on the different anatomical parts of a hand including metacarpophalangeal (MCP), proximal interphalangeal (PIP), distal interphalangeal (DIP), and fingertip (Fin).

**Results.** The quantitative results are presented in Table I. In comparison, MAPS achieves the highest average accuracy
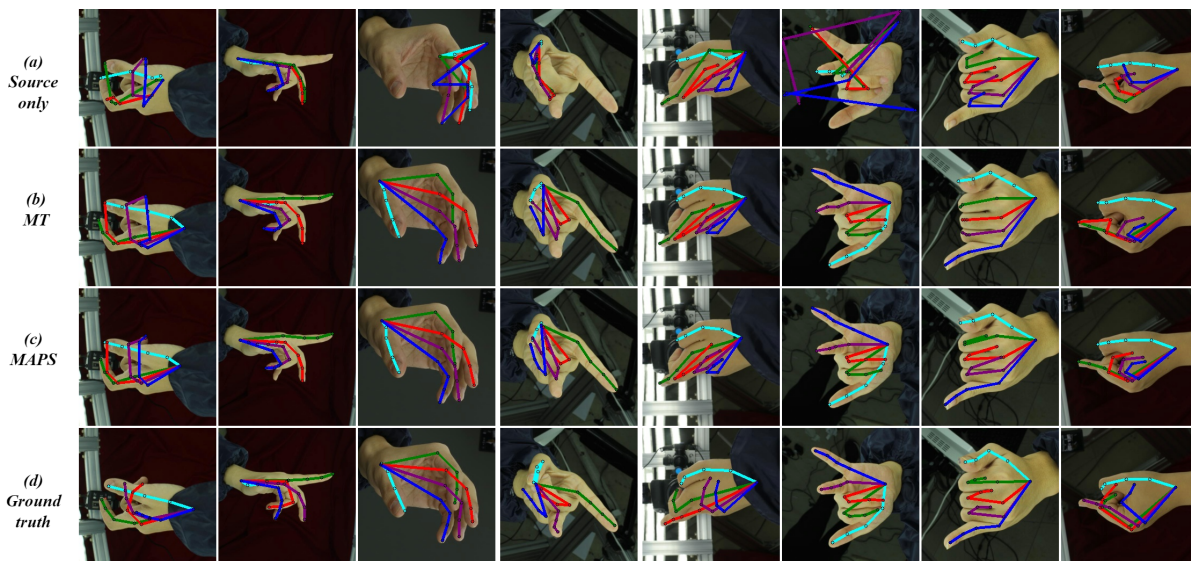
Fig. 2: Qualitative results of task **RHD→H3D**.

TABLE II: PCK@0.05 on task **SURREAL→LSP**.

| Method | SF | Sld | Elb | Wrist | Hip | Knee | Ankle | Avg. |
|---|---|---|---|---|---|---|---|---|
| CCSSL [54] | × | 36.8 | 66.3 | 63.9 | 59.6 | 67.3 | 70.4 | 60.7 |
| UDA-Animal [18] | × | 61.4 | 77.7 | 75.5 | 65.8 | 76.7 | 78.3 | 69.2 |
| RegDA [19] | × | 62.7 | 76.7 | 71.1 | 81.0 | 80.3 | 75.3 | 74.6 |
| UDAPE [20] | × | **69.2** | **84.9** | **83.3** | **85.5** | **84.7** | **84.3** | **82.0** |
| Source only | - | 50.6 | 64.8 | 63.3 | 70.1 | 71.2 | 70.1 | 65.0 |
| MT | ✓ | **69.2** | 82.1 | 80.2 | 83.4 | 82.3 | 80.8 | 79.6 |
| MAPS | ✓ | 67.0 | **84.2** | **82.7** | **84.0** | **83.8** | **84.6** | **81.0** |

TABLE III: PCK@0.05 on task **SynAnimal→AnimalPose**.

| Method | SF | Dog | | | | | Sheep | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Eye | Hoof | Knee | Elb | Avg. | Eye | Hoof | Knee | Elb | Avg. |
| CCSSL [54] | × | 34.7 | 37.4 | 25.4 | 19.6 | 27.0 | 44.3 | 55.4 | 43.5 | 28.5 | 42.8 |
| UDA-Animal [18] | × | 26.2 | 39.8 | 31.6 | 24.7 | 31.1 | 48.2 | 52.9 | 49.9 | 29.7 | 44.9 |
| RegDA [19] | × | 46.8 | 54.6 | 32.9 | 31.2 | 40.6 | 62.8 | 68.5 | 57.0 | 42.4 | 56.9 |
| UDAPE [20] | × | **56.1** | 59.2 | **38.9** | **32.7** | **45.4** | 61.6 | **77.4** | 57.7 | **44.6** | 60.2 |
| Source only | - | 38.2 | 43.2 | 25.7 | 24.1 | 32.0 | 59.9 | 60.7 | 46.2 | 31.0 | 47.9 |
| MT | ✓ | 53.0 | **59.3** | **37.7** | 32.0 | 44.4 | **66.1** | 76.5 | **59.8** | 42.6 | **60.6** |
| MAPS | ✓ | **63.2** | **60.5** | 37.3 | **34.0** | **46.7** | **74.2** | **77.9** | 58.8 | **43.3** | **62.0** |

even compared with the non-source-free methods, and MT also achieves competitive results. We outperform a large margin than most of the methods except UDAPE [20], demonstrating the effectiveness of our method. MAPS also outperforms MT, as the noise accumulation phenomenon in MT is reduced. We show the qualitative results in Fig. 2. The results of 'source-only' are not satisfied in some complicated hand poses due to a large domain gap, and MT successfully aligns the target domain to the unseen source domain. MAPS further refines the results of MT, obtaining more accurate results.

### C. Human Keypoint Detection

**Dataset.** SURREAL [70] is a synthetic dataset that consists of monocular videos of people with a total of 6M images,

which are photo-realistic renderings of people under large variations in shape, texture, viewpoint, and pose. Leeds Sports Pose [71] (**LSP**) is widely used as the benchmark for human keypoint detection. LSP contains a total of 2k images of sports persons with annotated human body joint locations gathered from Flickr.

**Implementation Details.** We use **SURREAL** as the source domain and **LSP** as the target domain. The increasing sequence $M$ in the baby step schedular is set to $\{0.1N, 0.2N, 0.3N\}$, where $N$ is the total number of training samples, and the number of progressive selection rounds $Q$ is 3. We report 16 keypoints on the human body parts, *i.e.*, shoulder (Sld), elbow (Elb), wrist, hip, knee, and ankle.

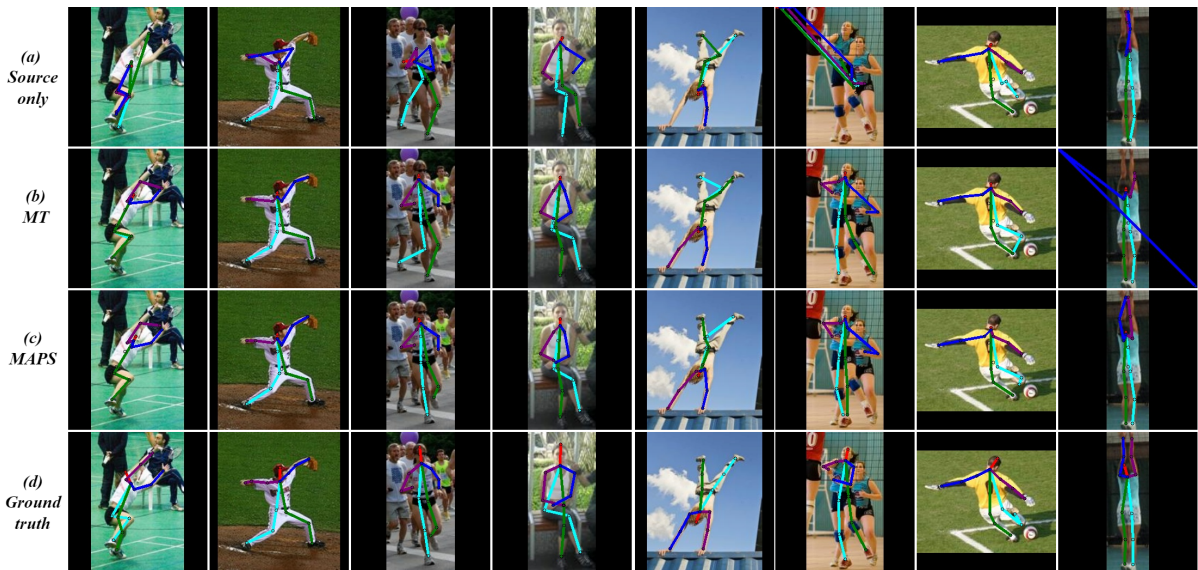**Results.** We show the quantitative results in Table II. In this

Fig. 3: Qualitative results of task **SURREAL→ LSP**.

TABLE IV: PCK@0.05 on task **SynAnimal→TigDog**.

| Method | SF | Horse | | | | | | | | Tiger | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Eye | Chin | Sld | Hip | Elb | Knee | Hoof | Avg. | Eye | Chin | Sld | Hip | Elb | Knee | Hoof | Avg. |
| CCSSL [54] | × | 89.3 | 92.6 | 69.5 | 78.1 | 70 | 73.1 | 65 | 73.1 | 94.3 | 91.3 | 49.5 | 70.2 | 53.9 | 59.1 | 70.2 | 66.7 |
| UDA-Animal [18] | × | 86.9 | 93.7 | 76.4 | 81.9 | 70.6 | 79.1 | 72.6 | 77.5 | 98.4 | 87.2 | 49.4 | 74.9 | 49.8 | 62 | 73.4 | 67.7 |
| RegDA [19] | × | 89.2 | 92.3 | 70.5 | 77.5 | 71.5 | 72.7 | 63.2 | 73.2 | 93.3 | 92.8 | 50.3 | 67.8 | 50.2 | 55.4 | 60.7 | 61.8 |
| UDAPE [20] | × | 91.3 | 92.5 | 74.0 | 74.2 | 75.8 | 77.0 | 66.6 | 76.4 | 98.5 | 96.9 | 56.2 | 63.7 | 52.3 | 62.8 | 72.8 | 67.9 |
| Source only | - | 82.0 | 90.0 | 59.2 | 79.5 | 65.8 | 66.9 | 57.7 | 67.4 | 85.4 | 81.8 | 44.6 | 70.8 | 39.6 | 48.4 | 55.5 | 54.8 |
| MT | ✓ | 80.9 | 90.7 | 71.3 | 75.5 | 72.9 | 74.0 | 66.1 | 73.5 | 98.7 | 91.5 | 47.9 | 56.5 | 48.2 | 60.5 | 68.3 | 63.9 |
| MAPS | ✓ | 82.3 | 91.4 | 74.1 | 76.5 | 71.8 | 74.2 | 65.3 | 73.7 | 99.5 | 93.8 | 55.1 | 60.7 | 48.7 | 60.9 | 71.4 | 66.0 |

TABLE V: The effectiveness of three loss functions and their interactions on task **SynAnimal→AnimalPose**.

| $\mathcal{L}_{mt}$ | $\mathcal{L}_{mix}$ | $\mathcal{L}_{spl}$ | SynAnimal→AnimalPose: Dog | | | | | SynAnimal→AnimalPose: Sheep | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Eye | Hoof | Knee | Elb | Avg. | Eye | Hoof | Knee | Elb | Avg. |
| | | | 38.2 | 43.2 | 25.7 | 24.1 | 32.0 | 59.9 | 60.7 | 46.2 | 31.0 | 47.9 |
| ✓ | | | 53.0 | 59.3 | 37.7 | 32.0 | 44.4 | 66.1 | 76.5 | 59.8 | 42.6 | 60.6 |
| ✓ | ✓ | | 56.6 | 59.5 | 39.2 | 33.2 | 45.8 | 69.5 | 78.2 | 59.3 | 42.6 | 61.4 |
| ✓ | | ✓ | 59.7 | 59.1 | 36.7 | 32.5 | 45.2 | 72.3 | 76.8 | 58.9 | 43.6 | 61.6 |
| ✓ | ✓ | ✓ | 63.2 | 60.5 | 37.3 | 34.0 | 46.7 | 74.2 | 77.9 | 58.8 | 43.3 | 62.0 |

task, MAPS does not perform as excellently as the hand keypoint detection benchmark but also obtains a competitive result. UDAPE achieves the highest accuracy, which benefits from its input-level and output-level alignment modules. In their input-level alignment module, the style and statistic information of source data is leveraged, which is unavailable in our setting. Besides, MAPS increases the accuracy of MT by 1.4%, which is an obvious improvement to validate the effectiveness of the proposed method. The qualitative results are shown in Fig. 3, we improve the ability of difficult points prediction such as knee and elbow.

### D. Animal Keypoint Detection

**Dataset.** Synthetic Animal Dataset [54] (**SynAnimal**) is a synthetic pose dataset rendered from CAD models, containing 5 different animals including the horse, tiger, sheep, hound,

and elephant. Each class contains 10k images. **AnimalPose** dataset [49] contains 6.1k images in the wild from 5 animals, *i.e.*, dog, cat, cow, sheep, and horse. **TigDog** dataset [72] is a large-scale animal keypoint detection dataset, which provides 30k images from real-world videos of horse and tigers.

**Implementation Details.** We conduct two experiments in this section, with **SynAnimal** as the source domain in both experiments, **AnimalPose** and **TigDog** as the target domain, respectively. In **SynAnimal→ AnimalPose**, the increasing sequence $M$ is set to $\{0.3N, 0.4N, 0.5N\}$, and the number of progressive selection rounds $Q$ is set to 3. We report the detection results of 14 keypoints on the dog and sheep body parts including the eye, hoof, knee, and elbow (Elb). In **SynAnimal→ TigDog**, $M$ is set to $\{0.1N, 0.15N, 0.2N\}$, and $Q$ is also set to 3. We present the detection results of 14 keypoints on the horse and tiger body parts including the eye,

TABLE VI: The effectiveness of three loss functions and their interactions on task **SURREAL→LSP**.

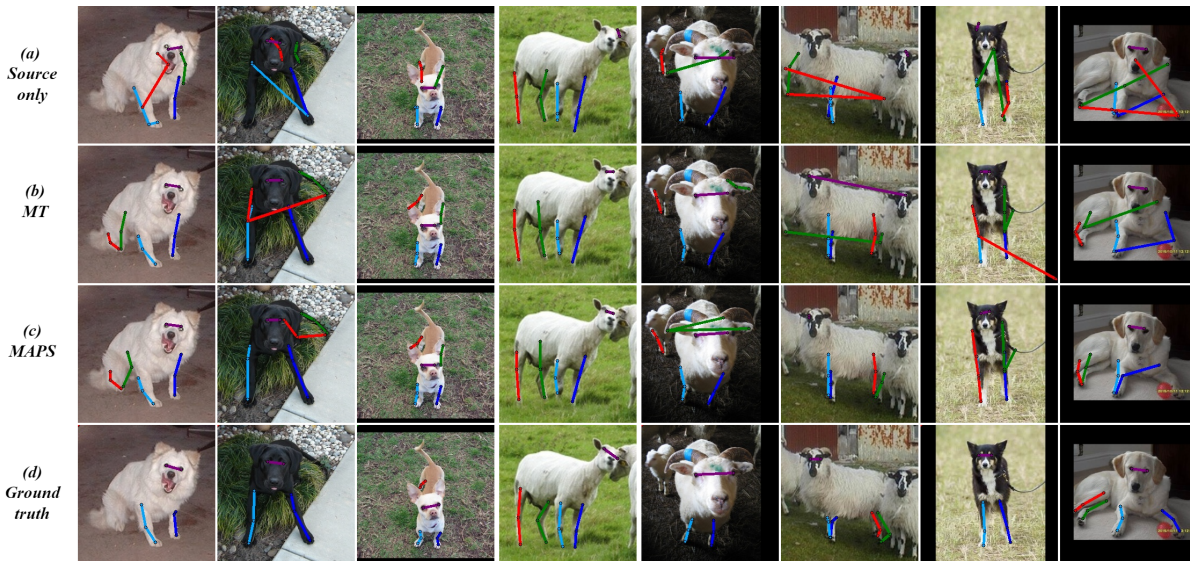| $\mathcal{L}_{mt}$ | $\mathcal{L}_{mix}$ | $\mathcal{L}_{spl}$ | SURREAL→LSP | | | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | | Sld | Elb | Wrist | Hip | Knee | Ankle | Avg. |
| | | | 50.6 | 64.8 | 63.3 | 70.1 | 71.2 | 70.01 | 65.0 |
| ✓ | | | **69.2** | 82.1 | 80.2 | 83.4 | 82.3 | 80.8 | 79.6 |
| ✓ | ✓ | | 69.1 | 83.7 | 77.7 | **85.2** | 84.3 | 82.5 | 80.4 |
| ✓ | | ✓ | 67.4 | 83.3 | 80.5 | 83.9 | 83.6 | 83.7 | 80.4 |
| ✓ | ✓ | ✓ | 67.0 | **84.2** | **82.7** | 84.0 | **83.8** | **84.6** | **81.0** |



Fig. 4: Qualitative results of task **SynAnimal→AnimalPose**.

TABLE VII: Ablation study on different sample selection strategies on task **SynAnimal→TigDog**. 'Full set': selecting the full target set as easy samples in one round. 'One round': selecting part of confident samples in one round. 'SPL': selecting part of samples by self-paced learning strategy of MAPS.

| Selection | SynAnimal→TigDog: Dog | | | | | SynAnimal→TigDog: Sheep | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | Eye | Hoof | Knee | Elb | Avg. | Eye | Hoof | Knee | Elb | Avg. |
| Full set | 58.5 | 60.7 | **39.4** | 32.5 | 46.2 | 69.5 | 76.5 | 55.6 | **46.6** | 61.0 |
| One round | 61.9 | **61.3** | 38.6 | **34.7** | **47.3** | 70.9 | 76.0 | 56.9 | 43.0 | 60.4 |
| SPL | **63.2** | 60.5 | 37.3 | 34.0 | 46.7 | **74.2** | **77.9** | **58.8** | 43.3 | **62.0** |

chin, shoulder (Sld), hip, elbow (Elb), knee, and hoof. We run these experiments with the same augmentation in UPADE.

**Results** The results of **SynAnimal→AnimalPose** are shown in Table III. MAPS achieves the best results among these methods on both dog and sheep keypoint detection and improves the baseline method MT by a large margin (2%). **AnimalPose** is a small dataset, which makes the model easily overfitting to outliers, and our self-mixup loss favors the predictions to vary linearly with the input to improve the robustness, therefore performing well on such a small dataset. The qualitative results are shown in Fig. 4. The results of **SynAnimal→ TigDog** are shown in Table IV. In this task, MAPS still has room for improvement. The reason why MAPS performs not well is probably due to the fact that the domain gap between **SynAnimal** and **TigDog** is primarily reflected in image style. UDAPE [20] directly solves the style alignment by introducing

a pretrained style transfer model. UDA-Animal [18], as a typical animal keypoint detection method, narrows the domain gap through adversarial training, which relies on the statistical information in source images. The lack of the style information of the source samples suppresses the performance of MAPS on this task, but MAPS still outperforms the MT, which demonstrates our method works in large-scale scenarios.
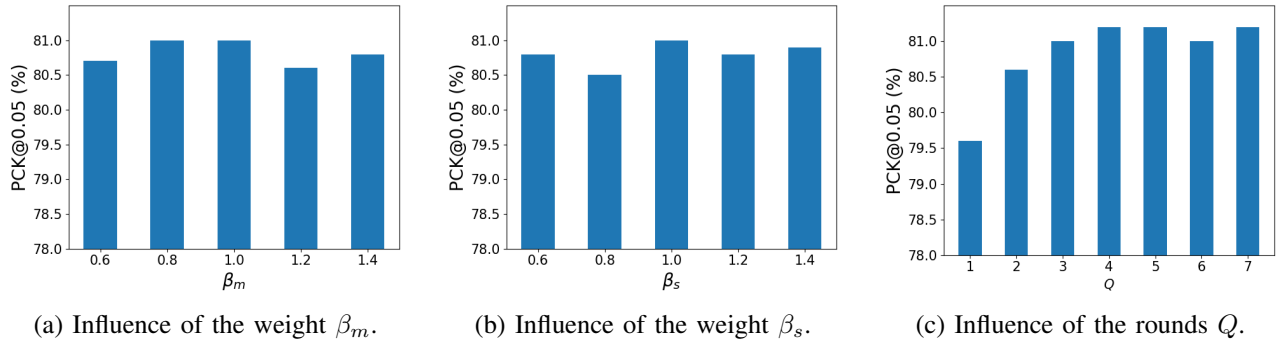
*E. Ablation Study*

To explore the effectiveness of the proposed strategies in different scenarios, For fairness, we conduct the experiments on two benchmarks, *i.e.*, **SynAnimal→AnimalPose** and **SURREAL→LSP**.

**Effectiveness of loss functions.** We first study the effectiveness of the three loss functions and the interactions among

TABLE VIII: Ablation study on different sample selection strategies on task **SURREAL→LSP**.

| Selection | SURREAL→LSP | | | | | | |
|---|---|---|---|---|---|---|---|
| | Sld | Elb | Wrist | Hip | Knee | Ankle | Avg. |
| Full set | 54.3 | 67.2 | 65.6 | 73.8 | 73.4 | 73.0 | 67.9 |
| One round | **65.3** | 81.1 | 79.3 | 83.8 | 84.0 | 84.1 | 79.6 |
| SPL | 64.4 | **84.3** | **81.5** | **84.4** | **84.5** | **84.5** | **80.6** |



(a) Influence of the weight $\beta_m$.

(b) Influence of the weight $\beta_s$.

(c) Influence of the rounds $Q$.

Fig. 5: Analysis of the weights $\beta_m$ and $\beta_s$, and the number of progressive selection rounds $Q$ on task **SURREAL→LSP**.

them. As shown in Table V and VI, the teacher-student consistency learning framework with the mean-teacher loss $\mathcal{L}_{mt}$ constructs a strong baseline method. The self-mixup loss $\mathcal{L}_m$ and the self-paced loss $\mathcal{L}_s$ improve the results of baseline by about 1% on the animal benchmark and 0.8% on the human benchmark, respectively. With both $\mathcal{L}_m$ and $\mathcal{L}_s$, the accuracy of the baseline is improved by about 2% on the animal benchmark and 1.6% on the human benchmark, demonstrating that both two proposed loss functions are practical, and the interactions among three loss functions are positive and useful.

**Effectiveness of sample selection strategies.** We conduct an ablation study on different sample selection methods. The results are presented in Table VII and VIII. In the tables, 'Full set' denotes selecting the full target set as easy samples in one round. 'One round' denotes selecting part of confident samples in one round. This can be seen as a special case of progressive selection, where $Q = 1$. 'SPL' denotes the progressive selection strategy in MAPS, *i.e.*, selecting part of samples by self-paced learning strategy. On both two tasks, the 'SPL' row obtains the highest accuracy, which demonstrates that selecting the part of confident samples in one round filters out most of the noisy pseudo labels, and the performance is improved. On task **SynAnimal→AnimalPose**, selecting the full target set as easy samples is better than selecting part of confident samples in one round. Although the full target set has more noisy labels, it also contains a larger amount of data and more information. The benefits of the large data and information can compensate for the negative effects of noise. On task **SURREAL→LSP**, the opposite is true. The negative effect of noise from using the full dataset is severe, and selecting part of confident samples alleviates this issue. In fact, this is a trade-off problem between noise and information, our progressive selection strategy improves the pseudo-label quality by leveraging the dynamic capacity of the model to achieve the best results.

### F. Further Analysis

**Sensitivity of the weights.** We conduct sensitivity analysis on two major hyper-parameters, *i.e.*, weights of loss functions, in our framework. We empirically choose value 1.0 for both two hyper-parameters. In Fig. 5 (a), the accuracy around $\beta_m = 1.0$ is stable at 80.5% to 81.0%. In Fig. 5 (b), the accuracy around $\beta_s = 1.0$ is also stable, floating around 0.5%. In summary, our weights of loss functions are not sensitive.

**Sensitivity of the number of selection rounds.** As shown in Fig. 5 (c), along with the increasing of progressive learning rounds $Q$, the accuracy increases from around 79.5% to around 81.5%. This is reasonable, because the smaller the rounds $Q$, the greater the noise of the selected easy samples, leading to worse performance. When it grows to $Q = 3$, the accuracy gradually stabilizes and does not bring greater gains as $Q$ grows. Therefore, we choose $Q = 3$ in our experiments.

## V. CONCLUSION

This paper considers a realistic setting called source-free domain adaptive keypoint detection, where only the well-trained source model is provided to the target domain. We first construct a simple baseline method based on the mean-teacher framework. Then we propose a new approach termed Mixup Augmentation and Progressive Selection (MAPS) built on this. The mixup augmentation regularizes the model to favor simple linear behavior in-between target samples thereby improving the robustness. The progressive selection strategy leverages the dynamic capacity of the current model to explore reliable pseudo labels. These two strategies respectively consider the convex behavior among samples and the quality of the pseudo label per sample, making them complement each other. Experiments verify that MAPS achieves competitive and even state-of-the-art performance in comparison to previous non-source-free counterparts.

## REFERENCES

[1] G. Csurka, "A comprehensive survey on domain adaptation for visual applications," *Domain adaptation in computer vision applications*, pp. 1–35, 2017. 1, 2

[2] M. Long, Y. Cao, J. Wang, and M. Jordan, "Learning transferable features with deep adaptation networks," in *Proc. ICML*, 2015. 1, 2

[3] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proc. CVPR*, 2017. 1, 2

[4] Z. Murez, S. Kolouri, D. Kriegman, R. Ramamoorthi, and K. Kim, "Image to image translation for domain adaptation," in *Proc. CVPR*, 2018. 1

[5] J. Liang, D. Hu, and J. Feng, "Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation," in *Proc. ICML*, 2020. 1, 2

[6] J. Huang, D. Guan, A. Xiao, and S. Lu, "Model adaptation: Historical contrastive learning for unsupervised domain adaptation without source data," *Proc. NeurIPS*, 2021. 1, 2

[7] S. Yang, J. van de Weijer, L. Herranz, S. Jui *et al.*, "Exploiting the intrinsic neighborhood structure for source-free domain adaptation," in *Proc. NeurIPS*, 2021. 1

[8] Y. Wang, J. Liang, and Z. Zhang, "Give me your trained model: Domain adaptive semantic segmentation without source data," *arXiv preprint arXiv:2106.11653*, 2021. 1, 2

[9] C. Yang, X. Guo, Z. Chen, and Y. Yuan, "Source free domain adaptation for medical image segmentation with fourier style mining," *Medical Image Analysis*, vol. 79, p. 102457, 2022. 1

[10] J. N. Kundu, A. R. Kulkarni, S. Bhambri, D. Mehta, S. A. Kulkarni, V. Jampani, and V. B. Radhakrishnan, "Balancing discriminability and transferability for source-free domain adaptation," in *Proc. ICML*, 2022. 1, 2

[11] X. Liu, C. Yoo, F. Xing, C.-C. J. Kuo, G. El Fakhri, J.-W. Kang, and J. Woo, "Unsupervised domain adaptation for segmentation with black-box source model," in *Medical Imaging 2022: Image Processing*, vol. 12032. SPIE, 2022, pp. 255–260. 1

[12] S. Li, M. Ye, X. Zhu, L. Zhou, and L. Xiong, "Source-free object detection by learning to overlook domain style," in *Proc. CVPR*, 2022. 1, 2

[13] X. Liu and Y. Yuan, "A source-free domain adaptive polyp detection framework with style diversification flow," *IEEE Transactions on Medical Imaging*, vol. 41, no. 7, pp. 1897–1908, 2022. 1

[14] D. Zhang, M. Ye, L. Xiong, S. Li, and X. Li, "Source-style transferred mean teacher for source-data free object detection," in *Proc. ACM-MMAsia*, 2021. 1

[15] Y. Ding, L. Sheng, J. Liang, A. Zheng, and R. He, "Proxymix: Proxy-based mixup training with label refinery for source-free domain adaptation," *arXiv preprint arXiv:2205.14566*, 2022. 1, 2

[16] J. Liu, X. Li, S. An, and Z. Chen, "Source-free unsupervised domain adaptation for blind image quality assessment," *arXiv preprint arXiv:2207.08124*, 2022. 1, 2

[17] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *Proc. CVPR*, 2017. 1

[18] C. Li and G. H. Lee, "From synthetic to real: Unsupervised domain adaptation for animal pose estimation," in *Proc. CVPR*, 2021. 1, 2, 4, 5, 6, 7, 8

[19] J. Jiang, Y. Ji, X. Wang, Y. Liu, J. Wang, and M. Long, "Regressive domain adaptation for unsupervised keypoint detection," in *Proc. CVPR*, 2021. 1, 2, 5, 6, 7

[20] D. Kim, K. Wang, K. Saenko, M. Betke, and S. Sclaroff, "A unified framework for domain adaptive pose estimation," in *Proc. ECCV*, 2022. 1, 2, 4, 5, 6, 7, 8

[21] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," *Proc. NeurIPS*, 2017. 1, 4

[22] D.-H. Lee *et al.*, "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks," in *Proc. ICML Workshops*, 2013. 1, 4

[23] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017. 1, 4

[24] M. Kumar, B. Packer, and D. Koller, "Self-paced learning for latent variable models," *Proc. NeurIPS*, 2010. 2, 3, 4

[25] S. Basu and J. Christensen, "Teaching classification boundaries to humans," in *Proc. AAAI*, 2013. 2

[26] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2009. 2

[27] M. Toldo, A. Maracani, U. Michieli, and P. Zanuttigh, "Unsupervised domain adaptation in semantic segmentation: a review," *Technologies*, vol. 8, no. 2, p. 35, 2020. 2

[28] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada, "Maximum classifier discrepancy for unsupervised domain adaptation," in *Proc. CVPR*, 2018. 2

[29] R. Li, Q. Jiao, W. Cao, H.-S. Wong, and S. Wu, "Model adaptation: Unsupervised domain adaptation without source data," in *Proc. CVPR*, 2020. 2

[30] J. Tian, J. Zhang, W. Li, and D. Xu, "Vdm-da: Virtual domain modeling for source data-free domain adaptation," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021. 2

[31] Z. Qiu, Y. Zhang, H. Lin, S. Niu, Y. Liu, Q. Du, and M. Tan, "Source-free domain adaptation via avatar prototype generation and adaptation," *arXiv preprint arXiv:2106.15326*, 2021. 2

[32] Y. Du, H. Yang, M. Chen, J. Jiang, H. Luo, and C. Wang, "Generation, augmentation, and alignment: A pseudo-source domain based method for source-free domain adaptation," *arXiv preprint arXiv:2109.04015*, 2021. 2

[33] H. Yan, Y. Guo, and C. Yang, "Source-free unsupervised domain adaptation with surrogate data generation," in *Proc. BMVC*, 2021. 2

[34] S. M. Ahmed, S. Lohit, K.-C. Peng, M. J. Jones, and A. K. Roy-Chowdhury, "Cross-modal knowledge transfer without task-relevant source data," *arXiv preprint arXiv:2209.04027*, 2022. 2

[35] N. Ding, Y. Xu, Y. Tang, C. Xu, Y. Wang, and D. Tao, "Source-free domain adaptation via distribution estimation," in *Proc. CVPR*, 2022. 2

[36] H. Yan, Y. Guo, and C. Yang, "Source-free unsupervised domain adaptation with surrogate data generation," in *Proc. BMVC*, 2021. 2

[37] F. Wang, Z. Han, Z. Zhang, and Y. Yin, "Active source free domain adaptation," *arXiv preprint arXiv:2205.10711*, 2022. 2

[38] J. Liang, D. Hu, J. Feng, and R. He, "Dine: Domain adaptation from single and multiple black-box predictors," in *Proc. CVPR*, 2022. 2

[39] M. Shen, Y. Bu, and G. Wornell, "On the benefits of selectivity in pseudo-labeling for unsupervised multi-source-free domain adaptation," *arXiv preprint arXiv:2202.00796*, 2022. 2

[40] N. Ma, J. Bu, J. Lu, J. Wen, Z. Zhang, S. Zhou, and X. Yan, "Semi-supervised hypothesis transfer for source-free domain adaptation," *arXiv preprint arXiv:2107.06735*, 2021. 2

[41] C. Eastwood, I. Mason, C. K. Williams, and B. Schölkopf, "Source-free adaptation to measurement shift via bottom-up feature restoration," *arXiv preprint arXiv:2107.05446*, 2021. 2

[42] J. Li, S. Bian, A. Zeng, C. Wang, B. Pang, W. Liu, and C. Lu, "Human pose regression with residual log-likelihood estimation," in *Proc. ICCV*, 2021. 2

[43] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *Proc. ECCV*, 2016. 2

[44] B. Xiao, H. Wu, and Y. Wei, "Simple baselines for human pose estimation and tracking," in *Proc. ECCV*, 2018. 2, 5

[45] H. Duan, K.-Y. Lin, S. Jin, W. Liu, C. Qian, and W. Ouyang, "Trb: a novel triplet representation for understanding 2d human body," in *Proc. ICCV*, 2019. 2

[46] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *Proc. CVPR*, 2019. 2

[47] T. Ohkawa, Y.-J. Li, Q. Fu, R. Furuta, K. M. Kitani, and Y. Sato, "Domain adaptive hand keypoint and pixel localization in the wild," *arXiv preprint arXiv:2203.08344*, 2022. 2

[48] Y. Zhang, T. Liu, M. Long, and M. Jordan, "Bridging theory and algorithm for domain adaptation," in *Proc. ICML*, 2019. 2

[49] J. Cao, H. Tang, H.-S. Fang, X. Shen, C. Lu, and Y.-W. Tai, "Cross-domain adaptation for animal pose estimation," in *Proc. ICCV*, 2019. 2, 7

[50] J. N. Kundu, S. Seth, P. YM, V. Jampani, A. Chakraborty, and R. V. Babu, "Uncertainty-aware adaptation for self-supervised 3d human pose estimation," in *Proc. CVPR*, 2022. 2

[51] C. Doersch and A. Zisserman, "Sim2real transfer learning for 3d human pose estimation: motion to the rescue," *Proc. NeurIPS*, vol. 32, 2019. 2

[52] T. Lee, B.-U. Lee, I. Shin, J. Choe, U. Shin, I. S. Kweon, and K.-J. Yoon, "Uda-cope: unsupervised domain adaptation for category-level object pose estimation," in *Proc. CVPR*, 2022. 2

[53] R. Jin, J. Zhang, J. Yang, and D. Tao, "Multi-branch adversarial regression for domain adaptative hand pose estimation," *IEEE Transactions on Circuits and Systems for Video Technology*, 2022. 2

[54] J. Mu, W. Qiu, G. D. Hager, and A. L. Yuille, "Learning from synthetic animals," in *Proc. CVPR*, 2020. 2, 5, 6, 7

[55] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proc. ICML*, 2009. 2

[56] Y. Fan, F. Tian, T. Qin, X.-Y. Li, and T.-Y. Liu, "Learning to teach," *arXiv preprint arXiv:1805.03643*, 2018. 2

[57] A. Graves, M. G. Bellemare, J. Menick, R. Munos, and K. Kavukcuoglu, "Automated curriculum learning for neural networks," in *Proc. ICML*, 2017. 2

[58] G. Hacohen and D. Weinshall, "On the power of curriculum learning in training deep networks," in *Proc. ICML*, 2019. 2

[59] X. Wang, Y. Chen, and W. Zhu, "A survey on curriculum learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 09, pp. 4555–4576, 2022. 2

[60] E. A. Platanios, O. Stretcu, G. Neubig, B. Poczos, and T. M. Mitchell, "Competence-based curriculum learning for neural machine translation," *arXiv preprint arXiv:1903.09848*, 2019. 3

[61] J. G. Tullis and A. S. Benjamin, "On the effectiveness of self-paced learning," *Journal of memory and language*, vol. 64, no. 2, pp. 109–118, 2011. 3

[62] D. Meng, Q. Zhao, and L. Jiang, "What objective does self-paced learning indeed optimize?" *arXiv preprint arXiv:1511.06049*, 2015. 3

[63] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler, "Joint training of a convolutional network and a graphical model for human pose estimation," *Proc. NeurIPS*, 2014. 4

[64] M. R. Elman, J. Minnier, X. Chang, and D. Choi, "Noise accumulation in high dimensional classification and total signal index." *Journal of Machine Learning Research*, vol. 21, pp. 36–1, 2020. 4

[65] H. Li and M. Gong, "Self-paced convolutional neural networks." in *Proc. IJCAI*, 2017. 4

[66] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016. 5

[67] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014. 5

[68] C. Zimmermann and T. Brox, "Learning to estimate 3d hand pose from single rgb images," in *Proc. ICCV*, 2017. 5

[69] Z. Zhao, T. Wang, S. Xia, and Y. Wang, "Hand-3d-studio: A new multi-view system for 3d hand reconstruction," in *Proc. ICASSP*, 2020. 5

[70] G. Varol, J. Romero, X. Martin, N. Mahmood, M. J. Black, I. Laptev, and C. Schmid, "Learning from synthetic humans," in *Proc. CVPR*, 2017. 6

[71] S. Johnson and M. Everingham, "Clustered pose and nonlinear appearance models for human pose estimation." in *Proc. BMVC*, 2010. 6

[72] L. Del Pero, S. Ricco, R. Sukthankar, and V. Ferrari, "Articulated motion discovery using pairs of trajectories," in *Proc. CVPR*, 2015. 7