

DOMAIN ADAPTIVE MULTIPLE INSTANCE LEARNING FOR INSTANCE-LEVEL PREDICTION OF PATHOLOGICAL IMAGES

Shusuke Takahama¹ Yusuke Kurose^{1,2} Yusuke Mukuta^{1,2} Hiroyuki Abe¹
 Akihiko Yoshizawa³ Tetsuo Ushiku¹ Masashi Fukayama⁴ Masanobu Kitagawa⁵
 Masaru Kitsuregawa^{1,6} Tatsuya Harada^{1,2,6}

¹ The University of Tokyo, Tokyo, Japan

² RIKEN Center for Advanced Intelligence Project, Tokyo, Japan

³ Kyoto University, Kyoto, Japan

⁴ The Japanese Society of Pathology, Tokyo, Japan

⁵ Tokyo Medical and Dental University, Tokyo, Japan

⁶ National Institute of Informatics, Tokyo, Japan

ABSTRACT

Pathological image analysis is an important process for detecting abnormalities such as cancer from cell images. However, since the image size is generally very large, the cost of providing detailed annotations is high, which makes it difficult to apply machine learning techniques. One way to improve the performance of identifying abnormalities while keeping the annotation cost low is to use only labels for each slide, or to use information from another dataset that has already been labeled. However, such weak supervisory information often does not provide sufficient performance. In this paper, we proposed a new task setting to improve the classification performance of the target dataset without increasing annotation costs. And to solve this problem, we propose a pipeline that uses multiple instance learning (MIL) and domain adaptation (DA) methods. Furthermore, in order to combine the supervisory information of both methods effectively, we propose a method to create pseudo-labels with high confidence. We conducted experiments on the pathological image dataset we created for this study and showed that the proposed method significantly improves the classification performance compared to existing methods.

Index Terms— Pathology, Multiple Instance Learning, Domain Adaptation

1. INTRODUCTION

In pathological diagnoses, doctors observe tissue slide images with a microscope and identify the presence of diseases such as cancer. Many studies have attempted to apply image recognition technology to reduce the burden on doctors through automatic diagnosis [1, 2]. Because the diagnosis requires detailed cell-level observation, the size of whole slide images (WSIs) can be as large as $10^5 \times 10^5$ pixels. Owing to memory

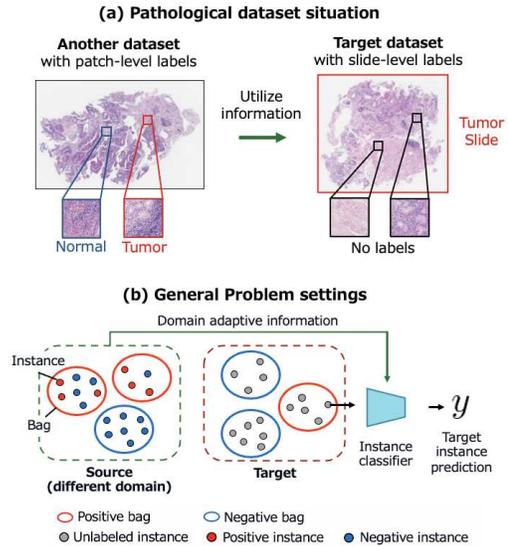


Fig. 1. (a) Our problem setting. We estimate the patch labels of the target dataset only with slide-level labels. We utilize information from another dataset that already has patch-level labels. (b) The generalized setting. A WSI can be represented as a bag, which is a set of instances. Our goal is to predict target instance labels by leveraging source information.

limitations, the WSI is often divided into small patch images to input classification models. Patch-level annotation takes a very high cost because it requires the expertise of doctors and a significant amount of time to annotate large WSIs. On the other hand, a label per slide, which indicates whether an abnormality exists in the WSI, requires little additional annotation cost. It is beneficial to improve the patch-level classification performance of the WSIs only with slide-level labels.

Multiple instance learning (MIL) is a type of weakly supervised learning with a single label for a bag of instances [3, 4, 5]. MIL methods have been applied to pathological im-

age analysis, regarding the patch image as the ‘‘instance’’ and the whole slide as the ‘‘bag’’ [6, 7]. Although this is very effective in reducing the annotation cost, the performance of the model trained only with slide labels was much lower than that with patch-level labels. On the other hand, using information from other datasets can also improve the classification performance without additional annotation costs. Domain adaptation (DA) is a method that utilizes a different domain to improve the performance of target data [8, 9, 10]. In pathological analysis, we can use existing public datasets with pixel-level labels such as the Camelyon dataset [1]. However, in most cases, we cannot use them directly because of the differences in body parts, appearance, and preprocesses such as tissue staining. Some studies have attempted to overcome the differences and transfer information between different pathological datasets [11, 12], but the performance is degraded when the difference between the domains is significant.

In this paper, we proposed a new problem setting to improve the patch-level classification performance of the target dataset only with slide labels, while utilizing information from the labeled source dataset from another domain (Fig. 1). However, since the supervised information from the source and target dataset is qualitatively different, there is no guarantee that simply combining the two will improve performance. Therefore, we propose a new training pipeline using pseudo-labels with high reliability by combining information from both the source and target. Our method can improve performance in situations where MIL alone and DA alone cannot provide accurate classification. We performed experiments on a new pathological dataset we created for this study, and the results confirmed that our method improves the instance classification performance compared to existing methods.

2. METHOD

Problem settings: In this study, we can access the labels of all instances from the source dataset, while we can only refer to bag labels and cannot access any instance labels of the target dataset. The purpose of this study is to estimate the instance labels of the target domain with high accuracy (Fig. 1 (b)). In the standard MIL setting, we consider bag $X = \{\mathbf{x}_1, \dots, \mathbf{x}_K\}$ as a set of instances $\mathbf{x} \in \mathbb{R}^D$. K is the number of instances in the bag, and it varies for each bag. Each instance has a binary label $y_k \in \{1, 0\}$, but this label cannot be referred to during training. The bag label $Y = 1$ if the bag contains at least one positive instance, and $Y = 0$ if instances are all negative. We can say that the source is a fully supervised setting, whereas the target is a standard MIL setting. We define the source domain $D_s = \{(X_i^s, Y_i^s)\}_{i=1}^{n_{sb}}$ as a set of n_{sb} bags, where the i -th bag $X_i^s = \{(\mathbf{x}_{ij}^s, y_{ij}^s)\}_{j=1}^{n_{si}}$ consists of n_{si} labeled instances, and the target domain $D_t = \{(X_i^t, Y_i^t)\}_{i=1}^{n_{tb}}$ as a set of n_{tb} bags, where the i -th bag $X_i^t = \{\mathbf{x}_{ij}^t\}_{j=1}^{n_{ti}}$ consists of n_{ti} unlabeled instances.

Overview of our method: Our pipeline consists of three

components: encoder G , bag classifier F_B , and instance classifier F_I . Each instance in the bag of the target and source \mathbf{x}_{ij} is input to G to obtain the feature vectors \mathbf{h}_{ij} . The feature vectors in the bag are collectively input to the bag classifier F_B to obtain the binary prediction score of the bag label $p(Y|X_i)$. By contrast, a feature vector from each instance is input to the instance classifier F_I to obtain the binary prediction score of the instance label $p(y|\mathbf{x}_{ij})$. Because the target does not have an instance label, the source instances with the instance labels and the target instance with the pseudo labels are used for training F_I . At the time of inference, we input the target instance features into F_I to obtain the prediction scores of the target instances $p(y|\mathbf{x}_{ij}^t)$.

$$\mathbf{h}_{ij} = G(\mathbf{x}_{ij}) \quad (1)$$

$$p(Y|X_i) = F_B(\mathbf{h}_{ij|j=1\dots n_i}) \quad (2)$$

$$p(y|\mathbf{x}_{ij}) = F_I(\mathbf{h}_{ij}) \quad (3)$$

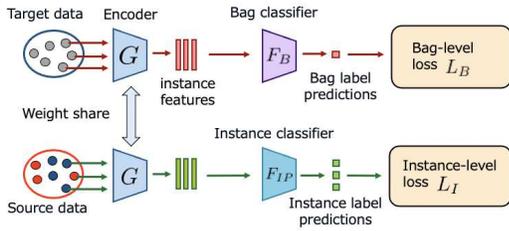
AttentionDeepMIL [5] is used as F_B . In this method, the bag feature is a weighted sum of instance features, and its weight a_{ij} is learnable. As mentioned in [5], the attention weight implies the positive score of each instance, so we used sigmoid instead of softmax to calculate a_{ij} so that we can directly obtain the positive score of an instance: $a_{ij} = \text{sigmoid}(\mathbf{w}^T \tanh(\mathbf{V}\mathbf{h}_{ij}^t))$. \mathbf{w} and \mathbf{V} are hyperparameters.

To improve the prediction performance of the target instances by F_I , we add a domain adaptation loss that performs distribution matching of the intermediate features \mathbf{h}_{ij} . We use MCD [9] as the DA loss. MCD performs feature distribution matching while considering category information by training the features to be away from the class boundary. To introduce MCD loss into our method, we use two instance classifiers F_{I1} and F_{I2} . We train G , F_{I1} and F_{I2} to minimize the instance classification loss $L_I(\mathbf{x}, y)$. At the same time, we train G to minimize the discrepancy loss $L_{adv}(\mathbf{x}^t)$ and two classifiers to maximize $L_{adv}(\mathbf{x}^t)$ alternately. The discrepancy loss is defined as $L_{adv}(\mathbf{x}^t) = \frac{1}{C} \sum_{i=1}^C |p_{1i} - p_{2i}|$, where $p_1(y|\mathbf{x}^t)$ and $p_2(y|\mathbf{x}^t)$ are the output of the two classifiers.

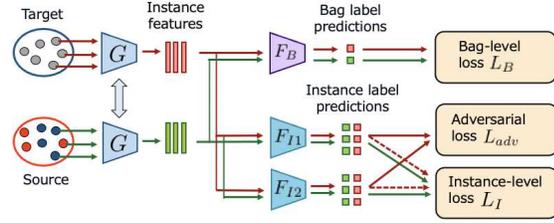
Pseudo labeling: Even if the feature distributions match, there will still be many misclassified instances if the decision boundaries of the source and target do not match. To tackle this problem, we directly optimize our model for target instance prediction by assigning pseudo-labels to the target instances and using them for the training of F_I . Because we know that all instances in the negative bag are negative, we mainly consider the instances from the positive bags. To obtain reliable pseudo-labels, we use two classifiers F_B and F_I . Because these two classifiers are trained using different supervisory information, they have different properties. we can obtain pseudo-labels with higher reliability by integrating the information from both of them.

We define the prediction score of the instance classifier $p_I(y|\mathbf{x})$ as the average of the predictions of F_{I1} and F_{I2} . We can also obtain the instance prediction score of the bag classifier $p_B(y|\mathbf{x})$ using the attention weight a_{ij} of F_B . Because

Step 1



Step 2



Step 3

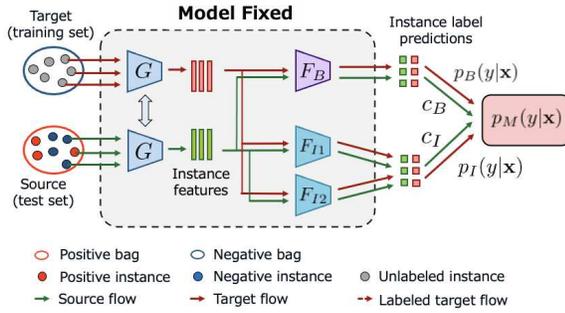


Fig. 2. The training pipeline. Step 1 is training with F_B and F_I separately. After converging Step 1, we alternately perform Steps 2 and 3. Step 2 involves training with DA loss and pseudo-labeled target instances. In Step 3, we assign pseudo-labels to target instances based on $p_M(y|\mathbf{x})$.

the two classifiers are trained in different ways, the accuracy of the predictions of both models can vary. For example, if the prediction performance of F_B is significantly poor, then the prediction of F_I should be mainly used. Therefore, we consider the confidence score of each model when assigning pseudo-labels. c_B and c_I represent the confidence scores of F_B and F_I , respectively. Because the target instances have no labels and we cannot directly examine the prediction accuracy, we instead use the PR-AUC score of the prediction performance of the source instances by each model as the confidence score. Then, we define the mix prediction score $p_M(y|\mathbf{x})$ as:

$$p_M(y|\mathbf{x}^t) = c_B * p_B(y|\mathbf{x}^t) + c_I * p_I(y|\mathbf{x}^t) \quad (4)$$

We first select positive candidates that satisfy $p_M(y = 0|\mathbf{x}^t) \leq p_M(y = 1|\mathbf{x}^t)$ and $0.5 \leq p_M(y = 1|\mathbf{x}^t)$, and negative candidates that satisfy $p_M(y = 1|\mathbf{x}^t) \leq p_M(y = 0|\mathbf{x}^t)$ and $0.5 \leq p_M(y = 0|\mathbf{x}^t)$. Then, from each of the positive and negative candidates, we select a fixed number of instances that

Table 1. The classification performance of each method on pathological dataset

	Accuracy	PR-AUC
Attention MIL	72.4±5.62	66.0±1.24
Source only	82.5±1.20	66.9±1.00
MCDDA	76.0±3.09	51.8±2.78
PLDA	78.6±1.66	76.5±3.97
Ours (Step 1)	82.7±2.84	71.1±2.62
Ours	86.0±4.11	83.4±3.48
Ideal case	91.1±0.79	87.1±2.05

have high confidence scores and assign pseudo-labels to them. We give pseudo-labels only to a small number of reliable instances at the beginning of the training when the prediction is ambiguous, and gradually increase the number as the training progresses. The definition of the number of pseudo-labels appears in the supplementary materials.

Training process: Figure 2 shows the entire pipeline of our method. Our method consists of three steps. In Step 1, we perform supervised learning of F_B using the target bag labels and F_I using the source instance labels. In this case, we use only one instance classifier F_{IP} . By performing Step 1, the training becomes more stable, and reliable pseudo-labels can be obtained from the beginning of Step 3. After Step 1 converges, we initialize two instance classifiers F_{I1} and F_{I2} and perform Steps 2 and 3 alternately. In Step 2, we train the model using the feature matching loss of DA. F_B is trained using both the source and target data. In addition, F_{I1} and F_{I2} are trained with $\{\mathbf{x}^m, \mathbf{y}^m\}$, which includes the source instances, the target instances with pseudo-labels from positive bags, and the sampled target instances from negative bags. We optimize the following three losses individually:

$$\min_{G, F_{I1}, F_{I2}, F_B} \lambda L_I(\mathbf{x}^m, \mathbf{y}^m) + (1 - \lambda)(L_B(X^t, Y^t) + L_B(X^s, Y^s)) \quad (5)$$

$$\min_{F_{I1}, F_{I2}} \lambda(L_I(\mathbf{x}^m, \mathbf{y}^m) - L_{adv}(\mathbf{x}^t)) \quad (6)$$

$$\min_G \lambda L_{adv}(\mathbf{x}^t) \quad (7)$$

where λ is a weight parameter. In Step 3, we fix the model parameters and give pseudo-labels to the target instances from positive bags. At the same time, we input the source instances in the validation set to calculate c_B and c_I in (4).

3. EXPERIMENT

In this section, we present the experimental results to confirm the effectiveness of our method. As a preliminary experiment, we performed detailed evaluations and ablation studies using benchmark datasets. The details are provided in the supplementary materials. In the following, we describe the results of the experiments using our pathological dataset.

Dataset: We constructed a new original dataset of pathological images to demonstrate the effectiveness of the proposed method. We collected whole slide images (WSI) from

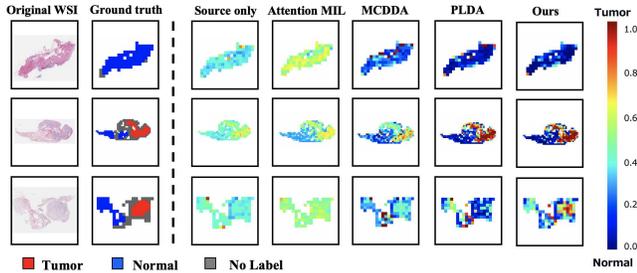


Fig. 3. Prediction heatmaps of positive prediction scores for target stomach dataset. In the ground-truth map, red indicates positive (anomaly), blue indicates negative (normal), and gray indicates areas without annotation.

two body parts, “Stomach” and “Colon,” which include 997 and 1368 WSIs, respectively. Many previous studies have set WSIs of two different datasets from a single organ as the source and the target, respectively [11, 12]. However, the domain gap between two organs is considerably larger than that in a single organ, making our settings more challenging and suitable for demonstrating the effectiveness of our method.

The size of the WSIs is approximately $10^4 \times 10^4$ pixels, and the maximum resolution is $\times 20$. Figure 1 (a) shows an example of the stomach WSI (right) and the colon WSI (left). The pixel-level normal/abnormal annotation was provided by expert pathologists. We separated WSI into patches of size 256×256 without overlaps and assign a binary label to each patch based on pixel-level annotation. The cropping, labeling, and image pre-processing methods followed the approach in [13]. We use Colon as the source and Stomach as the target.

Next, we created bags from each slide. Because one slide contains, at most, several hundred patches, we sampled 30 patches to make one bag. For the target positive slide, there’s no guarantee a bag includes positive instance because we don’t have patch labels. To increase the probability that the bag contains a positive instance, we performed clustering. First, we obtained the features of target patches using a classifier trained with source instances. Then, we separated the features into 10 clusters using K-means, and selected three samples with high positive scores from each cluster to obtain a bag of size 30. As a result, the probability that at least one positive instance is included in a bag created from the positive slide (confidence level of positive bag labels) was 95.0%, which is sufficiently reliable. For the negative slide, because all the patches were negative, we made as many bags as possible by randomly selecting patches. We randomly separated our dataset into 70% training slides and 30% test slides. Finally, we obtained 1000 training bags and 200 test bags from the source and target dataset, respectively.

Comparison methods: Since our problem setting is completely new, it cannot be compared to SOTA methods directly. To verify the effectiveness of MIL, DA, and pseudo-label modules, we evaluate the following comparison methods.

- **Attention MIL:** Train AttentionDeepMIL [5] only with

the target data. We used the values of the attention weights as the instance prediction scores.

- **Source only:** Train G and F_I with only source instances.
- **MCDDA:** Train unsupervised DA pipeline of MCD [9] without the bag labels of target data.
- **PLDA:** Train unsupervised DA with pseudo-labels. We assign pseudo-labels to the target instances by the model trained with the source and use them for the training.
- **Ours (Step 1):** Train only Step 1 of the proposed method. We evaluate the classification performance of F_{IP} .
- **Ideal case:** Train G and F_I using target instance labels that are not actually available. This is considered as the upper bound of the classification performance.

Experimental Settings: We used accuracy and the AUC of the precision-recall curve (PR-AUC) for instance-level prediction as the evaluation metrics. All experiments were conducted three times with random initial model weights, and the mean and standard deviation were calculated. We used ResNet50 [14] pretrained with ImageNet [15] as G , and two fully connected layers as F_I . The dimension of the output of G was 500. We trained 50 epochs for pretraining and source only, and 100 epochs for others. We set $\lambda = 0.5$.

Results: Table 1 shows the results of each method. Our proposed method outperforms other methods and achieves comparable scores with “Ideal case.” Figure 3 shows heatmaps of the estimated patch labels in the WSIs in the target test set by each trained model. The heatmaps of “Source only” and “Attention MIL” show little difference between the scores of the normal and abnormal areas, which implies that the predictions appear relatively vague. The heatmaps of “MCDDA” and “PLDA” appear to be relatively reasonable, but there are some regions of high abnormality scores in the normal region. This result is unfavorable for practical purposes because doctors need to examine the slide even if there is a small abnormality area. And “MCDDA” and “PLDA” do not detect the abnormal region well in the bottom example. Our proposed method made qualitatively valid prediction maps with a clear difference between the prediction scores of the normal and abnormal regions. Our method proved to be effective even in real-world applications such as pathological images.

4. CONCLUSION

In this study, we proposed a new problem setting to improve the classification performance of pathological images with low annotation cost, using only slide-level labels and information of another dataset from a different domain. In addition, we proposed a new pipeline to achieve the accurate classification of target instances by assigning pseudo-labels using two different supervisory information. Our method was evaluated on the pathological image dataset constructed in this study. The results demonstrate that our proposed method can achieve higher performance than comparative methods.

5. ACKNOWLEDGMENTS

This work was partially supported by AMED JP181k1010028 · JP191k1010036, JST AIP Acceleration Research JPMJCR20U3, Moonshot R&D Grant Number JPMJPS2011, CREST Grant Number JPMJCR2015, JSPS KAKENHI Grant Number JP19H01115 · JP19K20369 and Basic Research Grant (Super AI) of Institute for AI and Beyond of the University of Tokyo.

6. COMPLIANCE WITH ETHICAL STANDARDS

This study was performed in line with the principles of the Declaration of Helsinki. Approval was granted by the Ethics Committee of The University of Tokyo (9/1/2021, 21-222)

7. REFERENCES

- [1] Peter Bandi, Oscar Geessink, Quirine Manson, Marcorry Van Dijk, Maschenka Balkenhol, Meyke Hermsen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, Aoxiao Zhong, et al., “From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge,” *IEEE transactions on medical imaging*, vol. 38, no. 2, pp. 550–560, 2018.
- [2] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghahfoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I Sánchez, “A survey on deep learning in medical image analysis,” *Medical image analysis*, vol. 42, pp. 60–88, 2017.
- [3] Marc-André Carbonneau, Veronika Cheplygina, Eric Granger, and Ghyslain Gagnon, “Multiple instance learning: A survey of problem characteristics and applications,” *Pattern Recognition*, vol. 77, pp. 329–353, 2018.
- [4] Xinggong Wang, Yongluan Yan, Peng Tang, Xiang Bai, and Wenyu Liu, “Revisiting multiple instance neural networks,” *Pattern Recognition*, vol. 74, pp. 15–24, 2018.
- [5] Maximilian Ilse, Jakub Tomczak, and Max Welling, “Attention-based deep multiple instance learning,” in *International Conference on Machine Learning*, 2018, pp. 2127–2136.
- [6] Le Hou, Dimitris Samaras, Tahsin M Kurc, Yi Gao, James E Davis, and Joel H Saltz, “Patch-based convolutional neural network for whole slide tissue image classification,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2424–2433.
- [7] Noriaki Hashimoto, Daisuke Fukushima, Ryoichi Koga, Yusuke Takagi, Kaho Ko, Kei Kohno, Masato Nakaguro, Shigeo Nakamura, Hidekata Hontani, and Ichiro Takeuchi, “Multi-scale domain-adversarial multiple-instance cnn for cancer subtype classification with unannotated histopathological images,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3852–3861.
- [8] Yaroslav Ganin and Victor Lempitsky, “Unsupervised domain adaptation by backpropagation,” in *International conference on machine learning*, 2015, pp. 1180–1189.
- [9] 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*, 2018, pp. 3723–3732.
- [10] Chaoqi Chen, Weiping Xie, Wenbing Huang, Yu Rong, Xinghao Ding, Yue Huang, Tingyang Xu, and Junzhou Huang, “Progressive feature alignment for unsupervised domain adaptation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 627–636.
- [11] Yue Huang, Han Zheng, Chi Liu, Xinghao Ding, and Gustavo K Rohde, “Epithelium-stroma classification via convolutional neural networks and unsupervised domain adaptation in histopathological images,” *IEEE journal of biomedical and health informatics*, vol. 21, no. 6, pp. 1625–1632, 2017.
- [12] Jian Ren, Ilker Hacihaliloglu, Eric A Singer, David J Foran, and Xin Qi, “Adversarial domain adaptation for classification of prostate histopathology whole-slide images,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018, pp. 201–209.
- [13] Shusuke Takahama, Yusuke Kurose, Yusuke Mukuta, Hiroyuki Abe, Masashi Fukayama, Akihiko Yoshizawa, Masanobu Kitagawa, and Tatsuya Harada, “Multi-stage pathological image classification using semantic segmentation,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 10702–10711.
- [14] 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*, 2016, pp. 770–778.
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

SUPPLEMENTARY MATERIALS: DOMAIN ADAPTIVE MULTIPLE INSTANCE LEARNING FOR INSTANCE-LEVEL PREDICTION OF PATHOLOGICAL IMAGES

Shusuke Takahama¹ Yusuke Kurose^{1,2} Yusuke Mukuta^{1,2} Hiroyuki Abe¹
 Akihiko Yoshizawa³ Tetsuo Ushiku¹ Masashi Fukayama⁴ Masanobu Kitagawa⁵
 Masaru Kitsuregawa^{1,6} Tatsuya Harada^{1,2,6}

¹ The University of Tokyo, Tokyo, Japan

² RIKEN Center for Advanced Intelligence Project, Tokyo, Japan

³ Kyoto University, Kyoto, Japan

⁴ The Japanese Society of Pathology, Tokyo, Japan

⁵ Tokyo Medical and Dental University, Tokyo, Japan

⁶ National Institute of Informatics, Tokyo, Japan

1. METHOD

1.1. Number of pseudo-labels

We give pseudo-labels only to a small number of reliable instances at the beginning of the training when the prediction is ambiguous, and gradually increase the number as the training progresses. τ is a variable for the upper limit of the number of pseudo-labels:

$$\tau = \min \left(N_{min} + m * \frac{N_{max} - N_{min}}{M}, N_{max} \right) \quad (1)$$

where m denotes the training epoch, and τ increases depending on m . We assign labels to $A_p * \tau$ positive instances and $A_n * \tau$ negative instances. M , N_{min} , N_{max} , A_p and A_n are the hyperparameters. For the target negative bags, all instances are guaranteed to be negative. To maintain a balance between the number of labeled instances from the positive and negative bags, instances were randomly sampled from negative bags to equal the number of labeled instances from the positive bags.

We set $M = 20$, $A_p = A_n = 1$, $N_{min} = \frac{n_{tpi}}{10}$, and $N_{max} = \frac{n_{tpi}}{4}$, where n_{tpi} is the sum of the instances in all target positive bags in the experiment with the pathological dataset.

2. EXPERIMENTS ON DIGIT DATASETS

2.1. Details

In this supplementary material, we present the results of preliminary experiments to confirm the effectiveness of our method. We performed detailed evaluations using digit datasets (MNIST [1] and SVHN [2]) and VisDA dataset

[3], a benchmark datasets for image analysis and domain adaptation.

First, we describe the experiments using the digit datasets. We created a bag by collecting images from MNIST [1] and SVHN [2], as in [4]. We define “9” as a positive class and the others as a negative class. The bag label is positive when at least one “9” is in it, and otherwise negative. The bag size follows a normal distribution with a mean of 10 and a variance of 2, and the number of positive instances in a positive bag follows a normal distribution with a mean of 1 and a variance of 1. In both datasets, we made 2000 training bags and 500 test bags, including an equal number of positive and negative bags by random sampling.

We used an encoder with three convolution layers, followed by a fully connected layer as G . F_I consists of two fully connected layers. The size of the output of the encoder was 500. We use Adam [5] as the optimizer, and the learning rate was $1e-4$. We trained 50 epochs for pre-training and source only, and 100 epochs for others. We set $\lambda = 0.5$, $A_p = 1$, $A_n = 3$, $N_{min} = \frac{n_{tpi}}{30}$, and $N_{max} = \frac{n_{tpi}}{10}$, where n_{tpi} is the sum of the instances in all target positive bags. Other settings and model structures were the same as the experiment on the pathological dataset.

2.2. Result

Table 2 lists the prediction scores for each method. In the setting where MNIST is set as the target, our method shows high performance, but the improvement is negligible because “Attention MIL” and “PLDA” already yielded good predictions. The reasons for the poor “MCDDA” performance may be an imbalance in the number of positive and negative instances and the lack of pretraining with ImageNet [6] of the model. On the other hand, in the setting where SVHN is set as the target, the proposed method performs considerably better

Table 1. The classification performance on VisDA dataset.

	plane	bycl	bus	car	horse	knife	mcycl	person	plant	sktbrd	train	truck	mean
Attention MIL	95.1	79.6	81.4	5.5	5.0	89.3	6.4	5.6	89.5	68.2	6.3	6.8	44.9
Source only	77.4	50.9	46.9	24.0	60.5	25.2	59.6	35.6	59.1	29.9	39.3	22.6	44.3
MCDDA	91.4	75.1	75.5	43.5	89.7	24.2	80.4	43.8	80.3	44.2	76.7	41.7	63.9
PLDA	96.9	86.0	79.5	52.2	94.1	89.0	86.2	64.2	89.9	70.4	84.4	56.4	79.1
Ours (Step 1)	92.2	75.0	78.3	40.1	87.5	28.3	83.4	47.5	85.0	50.1	75.5	40.5	65.3
Ours	98.3	88.1	84.7	59.5	95.5	93.7	88.4	73.6	94.2	80.1	89.8	60.6	83.9
Ideal case	99.0	91.4	88.7	67.5	96.2	96.6	90.3	83.4	95.6	93.0	93.2	71.3	88.9

Table 2. The classification performance on the digit dataset.

Source	SVHN	MNIST
Target	MNIST	SVHN
Attention MIL	99.5±0.39	24.1±1.13
Source only	37.4±1.09	14.6±1.46
MCDDA	11.5±1.54	9.5±1.20
PLDA	98.0±1.81	53.6±1.51
Ours (Step 1)	97.1±0.56	30.6±2.27
Ours	99.4±0.56	80.7±0.56
Ideal case	99.8±0.00	86.7±0.29

than the comparative methods, because SVHN is more difficult to classify than MNIST. Our method achieves PR-AUC close to the “Ideal case,” owing to the complementary effect of the target bag-level loss and source instance-level loss. We found that our method achieved a high classification performance without additional annotation of the target instances.

3. EXPERIMENTS ON VISDA 2017

3.1. Details

In this section, we validate our proposed method using VisDA 2017 [3]. VisDA 2017 is one of the largest cross-domain datasets for object classification and consists of a training dataset of synthetic images, a validation dataset of real images collected from MSCOCO [7], and a test dataset of real images collected from a different domain from the validation dataset. Because the label of the test dataset is not available, we used the training dataset as the source domain and the validation dataset as the target domain. Because there are 12 categories, we evaluated the performance when one category was set as positive and the others as negative for every category.

In each setting, we created 500 training bags and 200 test bags. The bag size follows a normal distribution with a mean of 10 and a variance of 2. The number of positive instances in a positive bag follows a normal distribution with a mean of 2 and a variance of 2. The learning rate is $1e-6$ for methods including feature distribution matching loss and $1e-5$ for others. We trained 50 epochs for pretraining and source only, and 100 epochs for others. Other settings and model structures were the same as the experiment on the pathological dataset.

Table 3. The result of the ablation study on VisDA dataset.

	bus	car	knife	sktbrd
Ours (Step 1)	78.3	40.1	28.3	50.1
w/o pseudo-label	79.2	52.3	51.2	57.9
w/o feature matching	84.0	57.5	92.7	76.4
pseudo-label with F_I	84.2	59.1	93.3	78.9
pseudo-label with F_B	84.9	15.3	93.8	11.5
pseudo-label w/o conf. score	84.6	57.5	93.5	78.3
pseudo-label with PFAN	84.0	57.1	91.7	76.7
Ours	84.7	59.5	93.7	80.1

3.2. Results

Table 1 shows the performance of our proposed method and the comparison methods when each category is set as positive. Our method performed better than the comparison methods in all settings. In particular, our method improves the performance of categories such as “car” and “truck,” which are easily mistaken for other vehicle categories and show extremely low performance with “Attention MIL.” Further, our method shows comparable scores with the “Ideal case” when “plane” and “horse” are set as positive. In addition, “Ours (Step 1),” which simply combines DA and MIL, does not achieve high performance, which demonstrates the effectiveness of our method.

3.3. Feature distribution

Figure 1 presents the visualization result for the distribution of intermediate features after training each method when the “plant” is set as positive. In “Source only,” the decision boundary of the target negative and positive are ambiguous, and the distribution of the source and target are completely separated. Although the separation of the target positive and negative is slightly improved in “MCDDA” and “PLDA,” it is not sufficient for accurate classification. Moreover, in “Attention MIL,” although the decision boundary of the target is clearer, the target distribution is completely separated from the source, and there is no guarantee that the target data can be successfully classified by the decision boundary of F_I . In our proposed method, the decision boundary of the target is clear, and the distributions of the source and target are well-

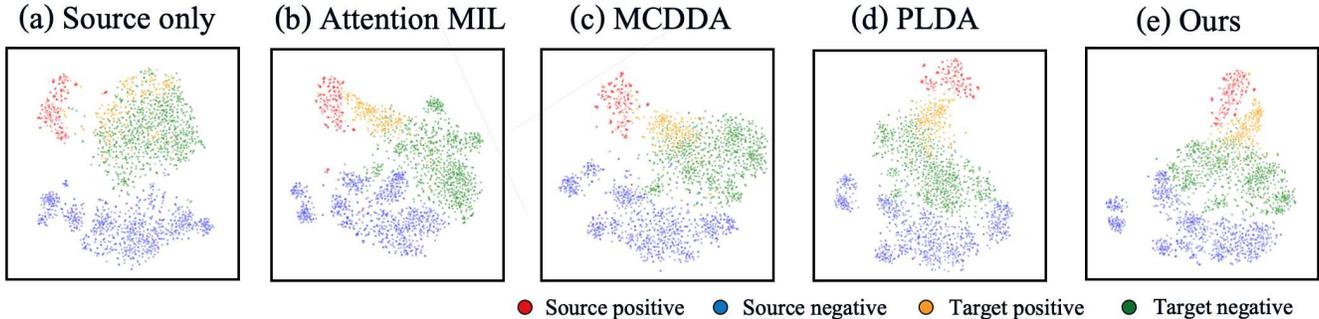


Fig. 1. Visualization of the feature of VisDA dataset when “plant” is positive. Red and blue indicate the source positive and negative instances, respectively, and yellow and green indicate the target positive and negative instances, respectively. We can observe that our proposed method achieves feature distribution matching and obtains a discriminative decision boundary.

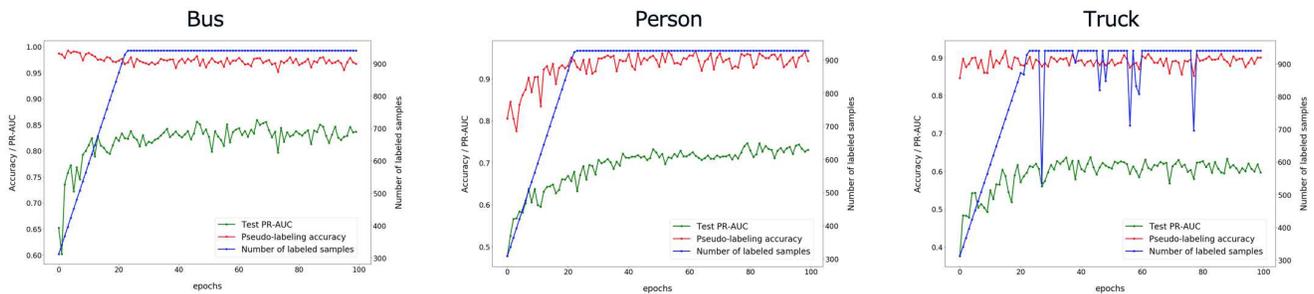


Fig. 2. Comparison of the number of labeled target instances, actual accuracy of pseudo-labels, and PR-AUC score during training on VisDA dataset.

matched. This means that F_I can achieve good classification performance, even in target instances.

3.4. Accuracy of pseudo-labels

Figure 2 shows the number of target instances given pseudo-labels, the actual accuracy of pseudo-labels, and the PR-AUC score of the proposed method for each training epoch on the VisDA dataset. The figure shows the results when “bus,” “person,” and “truck” are set as positive. The upper limit of the number of labeling instances was gradually increased up to the 20th epoch, and in most cases, the upper limit number of instances was labeled. Sometimes, as in the case of the “truck,” the number of labeled instances was lower than the upper limit, but the performance of the model improved steadily even in such cases. In the case of “bus,” although the pseudo-label accuracy decreased slightly as the number of labeled instances increased, it was maintained at a high level. By contrast, in the case of “person,” the pseudo-label accuracy increased as the model performance improved.

3.5. Ablation study

In our proposed method, the predictions of F_B and F_I are combined based on the confidence scores of the model to ob-

tain $p_M(y|\mathbf{x}^t)$, and pseudo-labels are assigned depending on the $p_M(y|\mathbf{x}^t)$ score. In addition, we added DA loss to achieve feature distribution matching. To clarify the effect of each module, we conducted ablation studies using the following methods:

- **w/o pseudo-label:** Train without pseudo-labels.
- **w/o feature matching:** Train without feature matching loss.
- **pseudo-label with F_I :** Assign pseudo-labels considering only the prediction of F_I .
- **pseudo-label with F_B :** Assign pseudo-labels considering only the prediction of F_B .
- **pseudo-label w/o conf. score:** Assign pseudo-labels without considering confidence scores c_I and c_B . We use the simple sum $p_M(y|\mathbf{x}^t) = p_B(y|\mathbf{x}^t) + p_I(y|\mathbf{x}^t)$.
- **pseudo-label with PFAN:** Giving pseudo-labels in feature space as in PFAN [8], which is one of the SOTA methods of pseudo-label DA. Specifically, we calculate the positive and negative centroids of the source instances and assign pseudo-labels to the target instances that are close to the centroids.

Table 3 shows the results of evaluating each method when “bus,” “car,” “knife,” and “skate” are set as positive. First, the performance of “w/o pseudo-label” and “w/o feature match-

ing” is degraded, indicating that pseudo-labels and feature matching loss contribute to the performance improvement. With respect to “pseudo-label with F_B ,” the performance for “bus” and “knife” is good, while the performance for “car” and “truck” is remarkably poor. This result is owing to the low confidence of F_B , which could not provide accurate pseudo-labels. By contrast, if the performance of F_I is extremely poor, then the performance of “the pseudo-label with F_I ” becomes worse. This indicates that the proposed method, which integrates the scores of both F_I and F_B considering the confidence score, contributes to performance improvement. In addition, the performance of “pseudo-label w/o conf. score,” which uses a simple sum of two predictions, is inferior to the proposed method. Furthermore, the proposed method outperforms PFAN which uses only a single score from feature space for labeling. This confirms the effectiveness of considering information from two models with different properties.

4. REFERENCES

- [1] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [2] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisaccho, Bo Wu, and Andrew Y Ng, “Reading digits in natural images with unsupervised feature learning,” in *In NIPS workshop on deep learning and unsupervised feature learning*, 2011, p. 5.
- [3] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko, “Visda: The visual domain adaptation challenge,” *arXiv preprint arXiv:1710.06924*, 2017.
- [4] Maximilian Ilse, Jakub Tomczak, and Max Welling, “Attention-based deep multiple instance learning,” in *International Conference on Machine Learning*, 2018, pp. 2127–2136.
- [5] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [7] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [8] Chaoqi Chen, Weiping Xie, Wenbing Huang, Yu Rong, Xinghao Ding, Yue Huang, Tingyang Xu, and Junzhou Huang, “Progressive feature alignment for unsupervised domain adaptation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 627–636.