# [Supplementary Material]
# H$^2$FA R-CNN: Holistic and Hierarchical Feature Alignment for Cross-domain Weakly Supervised Object Detection

Yunqiu Xu[1,2*]    Yifan Sun[1]    Zongxin Yang[3]    Jiaxu Miao[3]    Yi Yang[3]

[1]Baidu Research    [2]ReLER, AAII, University of Technology Sydney    [3]CCAI, Zhejiang University

imyunqiuxu@gmail.com  sunyf15@tsinghua.org.cn  {yangzongxin,jiaxumiao,yangyics}@zju.edu.cn

In this appendix, we provide details, results and discussions omitted in the main text:

- **Appendix A**: additional implementation details of 1) H$^2$FA R-CNN, 2) several baselines, and 3) our re-implemented CDWSOD method [9].

- **Appendix B**: additional experimental results, including 1) more detailed ablation study results, 2) more detailed comparison on similar domains, 3) more detailed benchmark results with different backbones, 4) loss weight sensitivity analysis, and 5) analysis on the selected regions by image-level recognition path.

- **Appendix C**: additional discussions, including: 1) the relationship between the image-level recognition path and WSDDN [2], 2) the intuition of the softmax along column, and 3) discussion on high-than-oracle results on several classes.

- **Appendix D**: additional visualizations, including 1) more detailed distribution visualization and 2) detection example visualization.

## A. Additional Implementation Details

**H$^2$FA R-CNN.** The pseudocode for the Instance- and Image-level Recognition (IIR) unit is shown in Algorithm 1. Since image-level annotations cannot provide effective spatial regularization for bounding-box regression, the regressors in RPN and detection head are trained with only source-domain data.

Our proposed method is implemented and evaluated using `Detectron2` [17] and `PaddleDetection` [1]. Widely used Faster R-CNN [10] with RoIAlign [6] is used as our base framwork. ImageNet [4] pre-trained ResNet-101 [7] is utilized as our network backbone in all experiments, unless otherwise specified. We use ResNet DC5 backbone variant by default, since of it achieves superior detection performance in most cases with faster training and inference speed, compared with the C4 counterpart. We use

---

*Work done during an internship at Baidu Research.

---

**Algorithm 1** Pseudocode of IIR unit, PyTorch-like

```
# feat: backbone features
# cls: classification logits
# obj: objectness logits
# target_domain: domain labels

obj, proposal = rpn(feat)
cls = det_head(feat, proposal)

if target_domain is True: # image-level recognition
    # exclude background class
    cls = cls[:, :-1]

    # objectness assignment
    idx = cls.argmax(dim=1)
    cls_obj = zeros_like(cls).scatter_(
        dim=1, index=idx, src=obj
    )

    # covert to image-level prediction, Eqn. (2)
    p = F_agg(cls, cls_obj)

    L = bce_loss(p, img_label)
else: # instance-level recognition
    L_rpn = bce_loss(sigmoid(obj), rpn_label)
    L_det = ce_loss(softmax(cls), det_label)
    L = L_rpn + L_det
```

**Notes**: The regression losses are omitted for clarity.

the VOC-standard AP (*i.e.*, IoU threshold is 50%) as evaluation metric.

Most training hyper-parameters are following the default configs of Faster R-CNN [10] in `Detectron2`. We use a mini-batch size of 8 (4 images per domain) in 2 NVIDIA V100 GPUs, and an initial learning rate is set to 0.005. The stochastic gradient descent (SGD) is utilized as the optimizer with a weight decay of 0.0001 and a momentum of 0.9.

For the Clipart$_{all}$, Clipart$_{test}$, Watercolor and Comic benchmarks [9], the `trainval` splits of PASCAL VOC 0712 [5] are used as source domain. Thus, we follow the default configs of `Detectron2` for training Faster R-CNN on PASCAL VOC. The image scale is [480, 800] pixels during training and 800 during inference. For the Clipart1k$_{all}$, we train for 36k iterations with the learning rate multiplied by 0.1 at 24k and 32k iterations. For the rest three data splits, we train for 24k iterations with the learning rate mul-

tiplied by 0.1 at 16k and 21.5k iterations.

For the noisy data setting, the `extra` splits of Watercolor and Comic provide additional ∼15.8k and ∼50.8k training images with noisy image-level annotations. Therefore, we enlarge the training iteration steps compared with original Watercolor and Comic benchmark. For the two noisy datasets, we train for 36k iterations with the learning rate multiplied by 0.1 at 24k and 32k iterations. Other hyper-parameters are the same as in original Watercolor and Comic benchmark.

For the similar domain adaptation setting (*i.e.*, from Cityscapes [3] to Foggy Cityscapes [12]), we follow the defaults configs of `Detectron2` for training Faster R-CNN on Cityscapes. We train with image scale (shorter side) randomly sampled from [800, 1024], which reduces overfitting; inference is on a single scale of 1024 pixels. We train for 24k iterations with the learning rate multiplied by 0.1 at 18k iterations.

**Source-only and oracle models.** The source-only and oracle models are trained with the default configs in `Detectron2`. The mini-batch size is set to 8 to fit in GPU memory. Accordingly, we adjust the initial learning rate to 0.01. The rest of the hyper-parameters (*e.g.*, optimizer, training iterations, training scales) are identical to that for H²FA R-CNN training. The source-only models are trained with instance-level labeled source domain. The oracle models are trained on the combination of instance-level labeled source and target domains, except for the oracle models for Foggy Cityscapes which are trained on only target domain as in previous UDAOD works.

**Reproducing related methods.** We re-implement a recent open-source CDWSOD method [9] in `Detectron2`. The re-implemented DT+PL [9] uses Faster R-CNN framework with ResNet-101 backbone as in our method. The mini-batch size is set to 8 to fit in GPU memory. The DT+PL is first trained on source domain for warming up, which shares the same training configs with the source-only models. As in [9], we then finetune detectors for one epoch and 10k iterations using the images obtained by DT (*i.e.*, domain transfer via a CycleGAN [19]) and PL (*i.e.*, pseudo label generation). The learning rate for detector finetuning is set to 0.0001.

# B. Additional Experimental Results

**Full results for ablation study.** Table B.1 summarizes the more detailed ablation study results. The top block reports the effectiveness of each feature alignment in H²FA R-CNN, when they are applied solely. Experimental results show all alignments improves the detection performance, except for A4 slight decrease the performance on Clipart$_{test}$.

| | A1 | A2 | A3 | A4 | Clipart$_{all}$ | Clipart$_{test}$ | Watercolor | Comic |
|---|---|---|---|---|---|---|---|---|
| (a) | | | | | 29.2 | 29.5 | 41.4 | 19.9 |
| (b) | ✓ | | | | 37.1 (+7.9) | 39.5 (+10.0) | 49.6 (+8.2) | 33.3 (+13.4) |
| (c) | | ✓ | | | 39.0 (+9.8) | 33.0 (+3.5) | 53.8 (+12.4) | 31.8 (+11.9) |
| (d) | | | ✓ | | 50.8 (+21.6) | 39.2 (+9.7) | 42.0 (+0.6) | 27.1 (+7.2) |
| (e) | | | | ✓ | 30.8 (+1.6) | 27.3 (-2.2) | 53.4 (+12.0) | 34.5 (+14.6) |
| (f) | ✓ | ✓ | | | 48.0 (+18.8) | 44.2 (+14.7) | 53.3 (+11.9) | 39.6 (+19.7) |
| (g) | ✓ | ✓ | ✓ | | 63.1 (+33.9) | 50.3 (+20.8) | 49.3 (+7.9) | 41.6 (+21.7) |
| (h) | ✓ | ✓ | | ✓ | 63.3 (+34.1) | 47.8 (+18.3) | 58.1 (+16.7) | 44.6 (+24.7) |
| (i) | | | ✓ | ✓ | 59.1 (+29.9) | 37.8 (+8.3) | 55.0 (+13.6) | 38.3 (+18.4) |
| (j) | | ✓ | ✓ | ✓ | 48.7 (+19.5) | 35.3 (+5.8) | 53.8 (+12.4) | 37.8 (+17.9) |
| (k) | ✓ | | ✓ | ✓ | 68.3 (+39.1) | 50.0 (+20.5) | 56.5 (+15.1) | 44.2 (+24.3) |
| (l) | ✓ | ✓ | ✓ | ✓ | 69.8 (+40.6) | 55.3 (+25.8) | 59.9 (+18.5) | 46.4 (+26.5) |

Table B.1. Effectiveness of different feature alignments in H²FA R-CNN, where mean AP performance (%) over all classes is reported. A1-A4 denote the four different types of feature alignments from bottom to top, where A1 and A2 are image-level alignments and A3 and A4 are instance-level alignments.

| | prsn | rider | car | truck | bus | train | mcycl | bcycl | mean |
|---|---|---|---|---|---|---|---|---|---|
| source-only | 34.9 | 41.8 | 44.1 | 14.1 | 27.3 | 9.1 | 31.6 | 41.4 | 30.5 |
| EPM [8] | 41.5 | 43.6 | 57.1 | 29.4 | 44.9 | **39.7** | 29.0 | 36.1 | 40.2 |
| FRCN w/ rot [14] | 45.8 | 51.0 | 63.1 | 26.8 | 47.1 | 23.6 | 30.6 | 43.6 | 41.5 |
| IIOD [15] | 32.8 | 44.4 | 49.6 | 33.0 | 46.1 | 38.0 | 29.9 | 35.3 | 38.6 |
| KTNet [13] | 43.0 | 42.7 | 60.0 | 32.3 | 46.6 | 38.4 | 31.2 | 38.2 | 41.5 |
| DT+PL [9] | 44.3 | **53.6** | 62.0 | 26.7 | 51.5 | 19.4 | 33.3 | 48.0 | 42.4 |
| H²FA R-CNN | **47.4** | 50.3 | **63.8** | **38.7** | **53.5** | 39.6 | **37.4** | **48.2** | **47.4** |
| oracle | 54.0 | 58.1 | 69.3 | 42.0 | 61.4 | 56.2 | 44.7 | 49.7 | 54.4 |

Table B.2. Mean AP performance (%) on Foggy Cityscapes with ResNet-101 backbone.

The second block summarizes the detection performance when coarse image-level alignments (A1 and A2) are imposed. After the coarse alignments are used, introducing fine-grained feature alignments (A3 or A4) consistently improves the detection performance. The third block shows the detection performance when fine instance-level alignments (A3 and A4) are applied. After introducing the most coarse-grained feature alignment (A1), detectors achieve significant improvement.

The bottom block shows the performance of H²FA R-CNN comprised of holistic feature alignments. When all alignments are imposed, extensive improvement is obtained. Remarkably, on Clipart$_{all}$, H²FA R-CNN boosts the source-only baseline from 29.2% mAP to 69.8%.

**Detailed results on similar domains.** The full results for similar domain adaptation (from Cityscapes to Foggy Cityscapes) are reported in Table B.2. The state-of-the-art UDAOD methods [8, 13–15] significantly improves source-only baseline by a large margin. The CDWSOD methods (*i.e.*, DT+PL [9] and H²FA R-CNN) surpass all previous UDAOD methods. H²FA R-CNN achieves the highest performance on 6 out of 8 classes, and obtains 47.4% overall mAP, exceeding the second place DT+PL [9] by 5.0%.

**Benchmark results with different backbones.** Table B.3 summarizes the comparison with five different backbones on five benchmarks. Concretely, we report the overall mAP of our H²FA R-CNN, source-only, and fully supervised oracle models. As shown in Table B.3, H²FA R-CNN

| backbone | VOC→Clipart$_{all}$ | | VOC→Clipart$_{test}$ | | | VOC→Watercolor | | | VOC→Comic | | | Cityscapes→Foggy Cityscapes | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | source | H$^2$FA R-CNN | source | H$^2$FA R-CNN | oracle | source | H$^2$FA R-CNN | oracle | source | H$^2$FA R-CNN | oracle | source | H$^2$FA R-CNN | oracle$^\dagger$ |
| VGG16 | 23.6 | 60.3 (+36.7) | 23.9 | 40.5 (+16.6) | 48.0 | 37.6 | 48.6 (+11.0) | 51.8 | 20.5 | 39.8 (+19.3) | 45.6 | 17.3 | 45.0 (+27.7) | 48.9 |
| R50-DC5 | 27.3 | 68.6 (+41.3) | 28.5 | 50.2 (+21.7) | 56.3 | 41.3 | 55.2 (+13.9) | 60.2 | 19.3 | 42.8 (+23.5) | 50.0 | 29.9 | 46.2 (+16.3) | 53.4 |
| R101-C4 | 27.0 | 66.6 (+39.6) | 27.7 | 51.2 (+23.5) | 58.5 | 42.8 | 57.2 (+14.4) | 60.1 | 20.1 | 41.5 (+21.4) | 57.6 | 26.2 | 48.9 (+22.7) | 54.4 |
| R101-DC5 | 29.2 | 69.8 (+40.6) | 29.5 | 55.6 (+26.1) | 59.3 | 41.4 | 59.9 (+18.5) | 59.9 | 19.9 | 46.4 (+26.5) | 53.7 | 30.5 | 47.4 (+16.9) | 54.4 |
| X101-DC5 | 25.9 | 73.6 (+47.7) | 27.1 | 53.9 (+26.8) | 59.6 | 45.2 | 59.8 (+14.6) | 60.6 | 19.7 | 46.6 (+26.9) | 54.2 | 30.2 | 49.2 (+19.0) | 54.8 |

Table B.3. Benchmark results with different backbones. Mean AP performance (%) over all classes is reported. † denotes oracle models are trained only on the instance-level labeled target domain.



(a) Loss weight for $\mathcal{L}_{ic}$



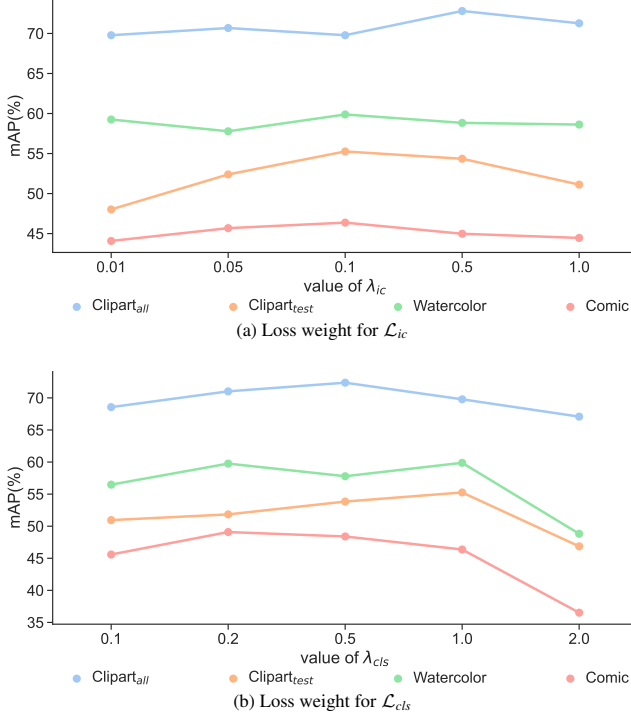(b) Loss weight for $\mathcal{L}_{cls}$

Figure B.1. Influence of different loss weights on four datasets, where mean AP performance (%) over all classes is reported.

achieves consistent and considerable performance improvement compared with source-only baseline. Compared with ResNet-101 C4 variant, ResNet-101 DC5 obtains higher performance on 4 out of 5 benchmarks. In addition, for detectors with the C4 variant, significantly more time consumption ($\sim 2\times$) is required for model training. We use the DC5-variant backbones in all other experiments.

**Loss weight sensitivity analysis.** There are five loss terms in the overall loss function $\mathcal{L}$ for an end-to-end optimization. We maintain the standard 1:1 weights for the detection losses (*i.e.*, $\mathcal{L}_{rpn}$ and $\mathcal{L}_{det}$) as in Faster R-CNN [10]. Following common practice [11, 16, 18], we set the weight for $\mathcal{L}_{dc}$ to 1. The 0.1 and 1 weights for $\mathcal{L}_{ic}$ and $\mathcal{L}_{cls}$ are empirically set. We analyze the impact of the last two weights, and find that our method has some robustness on their setting, as shown in Figure B.1. We use the same loss weights for all target domains in our experiments.

| | Clipart$_{all}$ | Clipart$_{test}$ | Watercolor | Comic |
|---|---|---|---|---|
| source-only | 45.4 | 46.6 | 58.8 | 31.3 |
| + img-level recog. path | 72.1 | 64.9 | 69.8 | **59.5** |
| H$^2$FA R-CNN | **79.5** | **80.1** | **73.1** | 58.1 |

Table B.4. CorLoc performance (%) on target-domain training data. After introducing instance-level alignments via image-level recognition path, CorLoc scores increase significantly, indicating the selected regions usually have high overlaps with ground-truths.

**Which regions are selected by the image-level recognition path.** In order to better understand how the image-level recognition path helps instance-level alignments, we analyze the correct localization (CorLoc) performance on target-domain training data. As summarized in Table B.4, image-level recognition path significantly improves CorLoc scores, indicating the selected regions usually have high overlaps with ground-truths.

Moreover, we visualize some heatmap examples of target-domain training data in Figure B.2. These heatmaps are calculated via accumulating the predicted scores (*i.e.*, objectness scores from RPN, classification scores from the detection head, and the aggregated scores in Figure 4) of each proposal and the corresponding proposal coordinates.

The class-agnostic objectness score maps (Figure B.2b) roughly highlight the foreground regions, while the classification score maps (Figure B.2c) coarsely locate some object regions in a class-aware manner. As shown in Figure B.2d, the highlighted regions in aggregated score maps are usually with less background noise and highly overlapped with corresponding foreground objects. Note that our image-level recognition path is not designed to mine all instances for a specific class from an input image. Instead, it considers representing the whole image with a few informative instances, as introduced in §3.3.

## C. Additional Discussions

**Relationship between the image-level recognition path and WSDDN [2].** There exists both resemblance and differences. The resemblance is that both approaches utilize image-level supervision for detection. Meanwhile, their differences are significant. Concretely, WSDDN [2] transforms the detection head to receive weak supervision and is not applicable to *RPN+detection head*. In contrast, our method re-uses the entire *RPN+detection head* to make image-level prediction. It allows simultaneous alignments

(a) Input images      (b) Objectness score maps      (c) Classification score maps      (d) Aggregated score maps

Figure B.2. Heatmap visualizations of the accumulated objectness scores, classification scores and aggregated scores for all proposals.

on RPN and detection head, which are critical to our method (see Table B.1).

**The intuition of the softmax operations in Figure 4 and Eqn. (2).** We use the two softmax operation for aggregating instance-level predictions into image-level predictions. Similar to the classifier in the detection head of Faster R-CNN [10], we first utilize a softmax along classes (*i.e.*, softmax along row in Figure 4, and $\sigma^{row}(\cdot)$ in Eqn. (2)) to extract the proposals' probabilities of belonging to each object class. Note that the *background* class is excluded during performing the softmax, as it does not explicitly exist in the image-level supervision.

Then, for each individual class, we use *weighted sum* to collect multiple proposal-level probability scores (*i.e.*, all the proposals' probabilities of belonging to this class) into a single image-level probability. To this end, we introduce a softmax across the proposals (*i.e.*, softmax along column in Figure 4, and $\sigma^{col}(\cdot)$ in Eqn. (2)) to generate the weight for each proposal. For the purpose of assigning weights, introducing a softmax is a common good choice, which provides normalization effect and naturally highlights the most representative proposals. In our preliminary experiments, removing the softmax makes the training fail to converge. The input of the softmax $\sigma^{col}(\cdot)$ (*i.e.*, the class-wise objectness $\bar{o}$) is derived from the objectness $o$ (*i.e.*, the output logits of the RPN). We empirically use 0 to initialize $\bar{o}$, as it shows better results than using a large negative number (*e.g.*, -10).

**Higher-than-oracle results on several classes.** We find that our method as well as some state-of-the-art methods achieve higher performance on several classes compared the oracle models. We conjecture it is due to the class imbalance problem on the target-domain datasets.

These tail classes have much fewer training instances (*e.g.*, 13 `bus` *vs*. 619 `person` instances on Clipart$_{test}$) and are easily wrongly recognized by the oracle. In contrast, our method partly alleviates this imbalance because the image-level annotations are less imbalanced (*e.g.*, 12 `bus` *vs*. 266 `person` images on Clipart$_{test}$), therefore improving the accuracy on tail classes.

Moreover, we also notice that on these tail classes, the performance gap looks large because the testing datasets contain very few samples *e.g.*, 8 *bus* instances on Clipart$_{test}$), as well. Successfully detecting 1 more instance could bring non-trivial improvement.

## D. Additional Visualizations

**Distribution visualization.** We extract the features of ground-truth bounding boxes and visualize the within-class distributions of different domains on Watercolor dataset. Figure D.3a shows source- and target-domain distributions of source-only baseline, in which the distributions of all classes are separated.

Figures D.3b-D.3e visualizes the distributions of four different feature alignments. When utilizing them separately, the distributions of two domains get closer, except for instance-level class-wise alignment. This indicates the importance of hierarchical feature alignment, *i.e.*, only perform fine-grained alignment at top could be unstable.

When all feature alignments are imposed, cross-domain features are aligned in a holistic and hierarchical manner. Consequently, H$^2$FA R-CNN achieves better aligning effect (see Figure D.3f). We also find that perfectly aligning cross-domain features remains difficult for some classes (*e.g.*, cat), even though the instance-level annotations of both domains are available (*i.e.*, the oracle model).

**Detection examples.** In Figures D.4-D.7, we visualize some detection results of our proposed H$^2$FA R-CNN on four benchmarks. Detections with confidence scores higher than 0.5 are visualized. H$^2$FA R-CNN is able to detect different instances in complicated scenes. H$^2$FA R-CNN can correctly localize more foreground regions, while other methods sometimes are confused due to severe domain shifts (see the third line in Figure D.6). Moreover, H$^2$FA R-CNN also has strong ability for distinguishing different different categories (see the third line in Figure D.7).

Figure D.8 shows several typical failure cases of H$^2$FA R-CNN. Without instance-level supervision guidance of target domain, it is difficult to reduce some false positive detections. For instance, H$^2$FA R-CNN sometimes gives some redundant predictions which cover small object parts or cover multiple objects. Small objects and crowded scenes remain challenging for H$^2$FA R-CNN.

## References

[1] PaddlePaddle Authors. PaddleDetection, object detection and instance segmentation toolkit based on PaddlePaddle. https://github.com/PaddlePaddle/PaddleDetection, 2019.

[2] Hakan Bilen and Andrea Vedaldi. Weakly supervised deep detection networks. In *CVPR*, 2016.

[3] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.

[4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.

[5] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (VOC) challenge. *IJCV*, 2010.

[6] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017.

[7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

Figure D.3. Within-class distributions of all classes on Watercolor [9] dataset, where the X-axis of each figure indicates the Euclidean distance to the center of source domain.
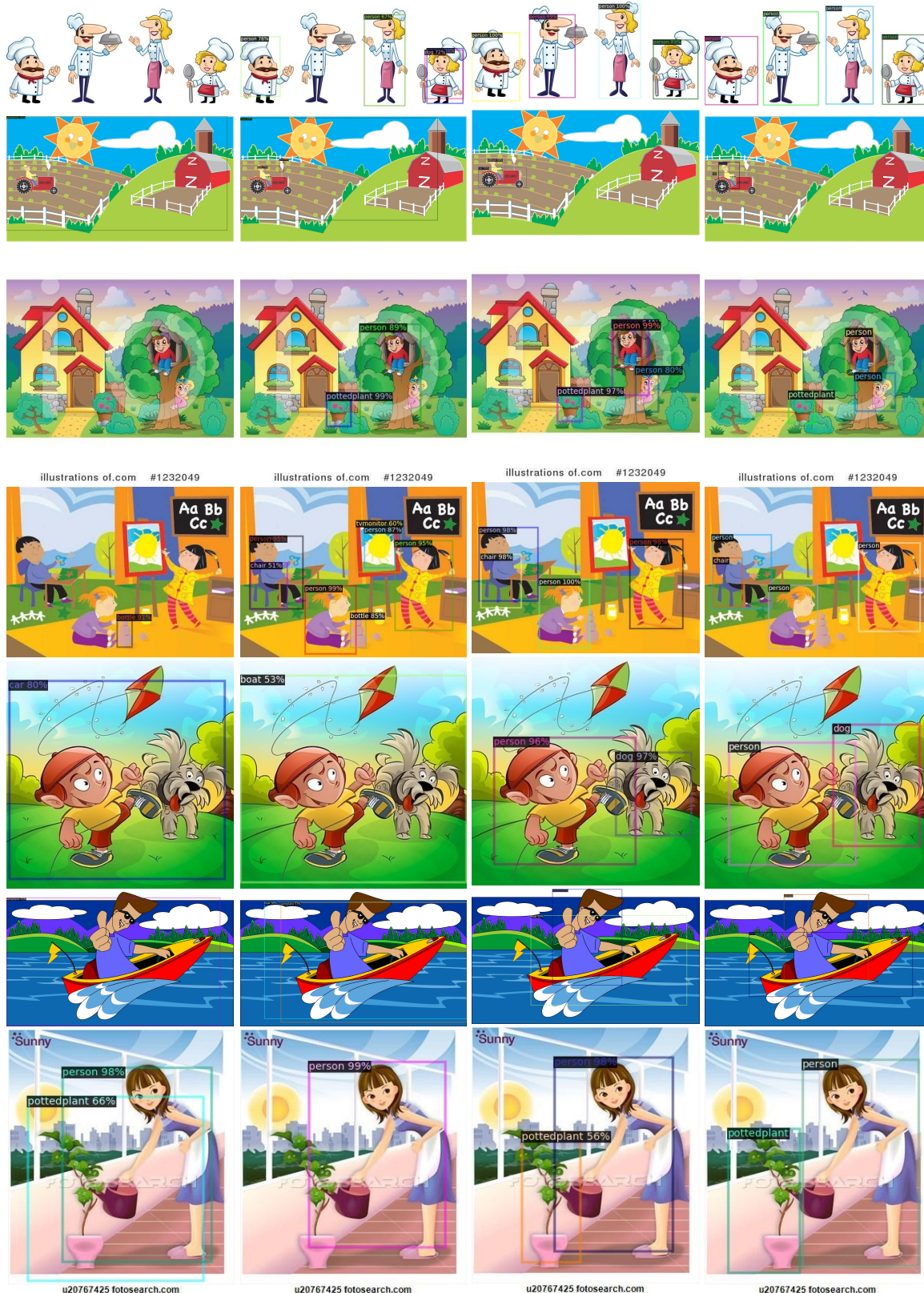
[8] Cheng-Chun Hsu, Yi-Hsuan Tsai, Yen-Yu Lin, and Ming-Hsuan Yang. Every pixel matters: Center-aware feature alignment for domain adaptive object detector. In *ECCV*, 2020.

[9] Naoto Inoue, Ryosuke Furuta, Toshihiko Yamasaki, and Kiyoharu Aizawa. Cross-domain weakly-supervised object de-

tection through progressive domain adaptation. In *CVPR*, 2018.

[10] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *NeurIPS*, 2015.

[11] Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate

(a) Source-only       (b) DT+PL [9]       (c) H$^2$FA R-CNN       (d) Ground-truth

Figure D.4. Detection examples on Clipart$_{all}$.
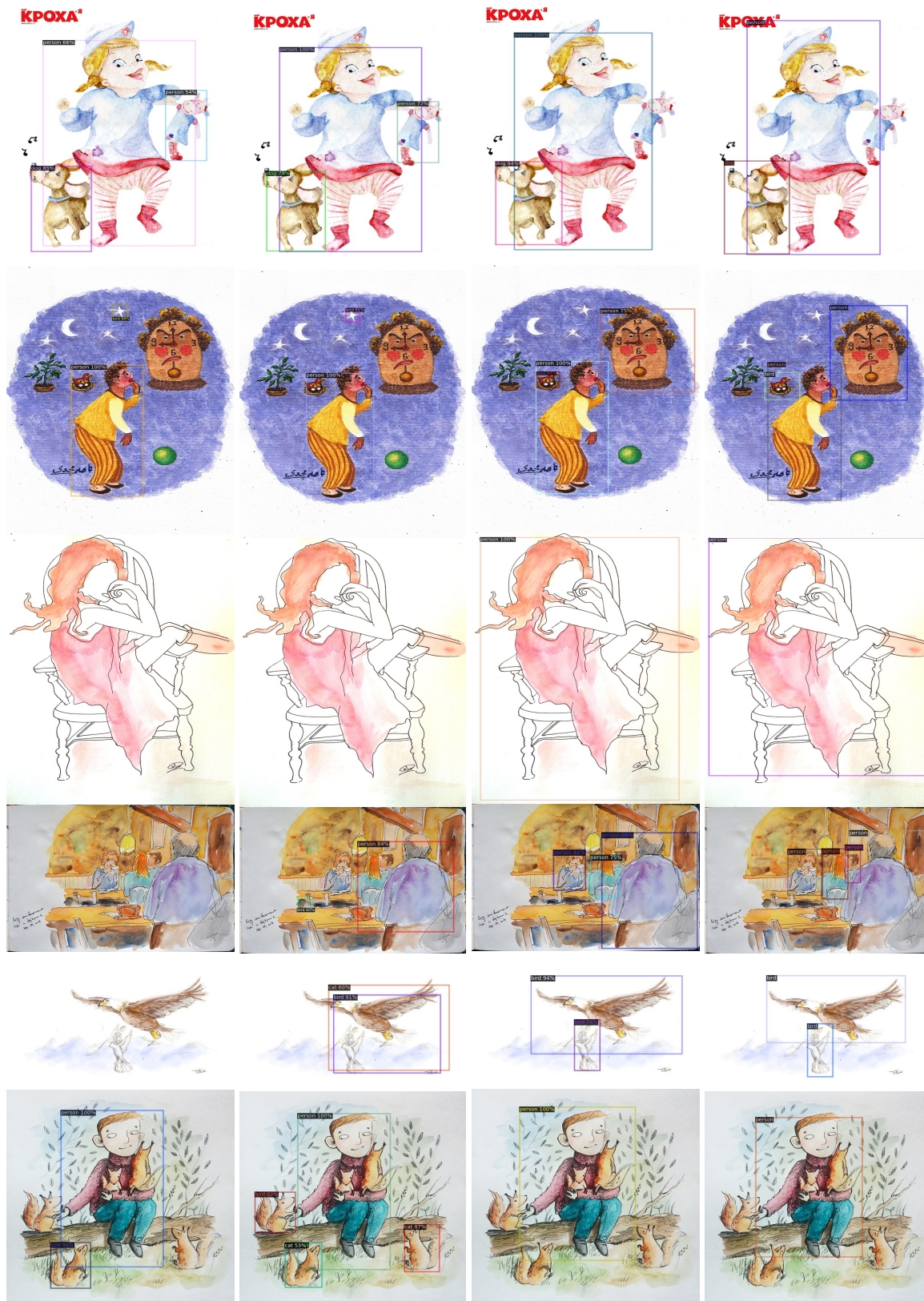
(a) Source-only      (b) DT+PL [9]      (c) H²FA R-CNN      (d) Ground-truth

Figure D.5. Detection examples on Clipart_test.

(a) Source-only      (b) DT+PL [9]      (c) H²FA R-CNN      (d) Ground-truth

Figure D.6. Detection examples on Watercolor.

(a) Source-only      (b) DT+PL [9]      (c) H²FA R-CNN      (d) Ground-truth

Figure D.7. Detection examples on Comic.

Figure D.8. Failure cases of H$^2$FA R-CNN. The ground-truth bounding-boxes are shown on the right.

Saenko. Strong-weak distribution alignment for adaptive object detection. In *CVPR*, 2019.

[12] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Semantic foggy scene understanding with synthetic data. *IJCV*, 2018.

[13] Kun Tian, Chenghao Zhang, Ying Wang, Shiming Xiang, and Chunhong Pan. Knowledge mining and transferring for domain adaptive object detection. In *ICCV*, 2021.

[14] Xin Wang, Thomas E Huang, Benlin Liu, Fisher Yu, Xiaolong Wang, Joseph E Gonzalez, and Trevor Darrell. Robust object detection via instance-level temporal cycle confusion. In *ICCV*, 2021.

[15] Aming Wu, Yahong Han, Linchao Zhu, and Yi Yang. Instance-invariant domain adaptive object detection via progressive disentanglement. *IEEE TPAMI*, 2021.

[16] Aming Wu, Rui Liu, Yahong Han, Linchao Zhu, and Yi Yang. Vector-decomposed disentanglement for domain-invariant object detection. In *ICCV*, 2021.

[17] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. https://github.com/facebookresearch/detectron2, 2019.

[18] Chang-Dong Xu, Xing-Ran Zhao, Xin Jin, and Xiu-Shen Wei. Exploring categorical regularization for domain adaptive object detection. In *CVPR*, 2020.

[19] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017.