

# Unsupervised Domain Adaptation for RGB-D Object Recognition

Matteo Bunino\*, Christian Cancedda† and Giacomo Garrone‡

\* matteo.bunino@studenti.polito.it

† christian.cancedda@studenti.polito.it

‡ giacomo.garrone@studenti.polito.it

**Abstract**—Convolutional neural networks require enormous amounts of data in order to achieve state of the art performances, however the data labeling process is a considerably expensive operation, well known to be one of the main bottlenecks in supervised learning. This issue can be partially addressed by relying on synthetic generated data, whose labeling is obtained without supervision, having the drawback of being slightly different from real data.

Domain adaptation techniques allow to learn robust domain invariant features from synthetic data during training by reducing the domain shift between synthetic and real data distributions.

This solution can be as well applied to the robotic vision field in which images come in both RGB and depth modalities, requiring the design of *ad-hoc* inter-modal self supervised tasks. In this paper we propose a novel self supervised task, comparing it to the model presented in [1], which is our main reference, and to the source-only baselines to prove its effectiveness versus simple training without domain adaptation.

We also propose a different initialization schema from [1] that improves its accuracy by 1.32% and try different approaches to enhance the performances of models, such as a data augmentation pipeline tailored to RGB-D data and a multi-task learning approach.

## I. INTRODUCTION

Convolutional neural networks rely on large amounts of image data in order to be trained properly. However, gathering and labeling real world images is a costly operation which is partially solved by synthetic generated data. This approach has the benefit that the examples can be labeled automatically without human supervision, but on the other hand it has the drawback that synthetic images are slightly different from their real counterpart.

Synthetic and real domains are said to be shifted among each other. Under the assumption that the two distributions from which data are generated only differ for a peculiar domain shift, domain adaptation techniques allow for the reduction of the effect of this factor. This was proposed for the first time in [2].

Robotics systems strongly rely on recognition tasks while interacting with the real world, therefore an effective solution that has been worked out is the use of RGB-D cameras that provide an additional information with respect to common cameras: the depth image. Depth images have a different nature from RGB ones: each pixel represents the distance from the camera, having a value which is not bounded in [0, 255] such as RGB pixels. The addition of depth images to the analyses allows us to take into consideration more

information than with RGB images alone. The latter have a rich content of texture details, whereas depth images convey more information about geometry and shape of the object. It is important to deploy a suited approach to depth images, due to their difference with RGBs, in fact in [3] it was proved that a CNN trained from scratch on depth images can learn very different features from the ones learned on RGB images.

Our work is focused on Domain Adaptation techniques applied to RGB-D images, aimed to reduce the domain shift between source (synthetic) and target (real) domains.

The current state of art is rich of proposed models that already employ domain adaptation, but most of them address this topic using RGB images only. In this case we want to tailor a domain adaptation model to the multimodal nature of RGB-D images. This requires to focus on the meaning of multimodal representation of an object.

A multimodal representation is the set of representations in different domains of the same object, that in this case it is composed by

- its visual appearance, that mapped on a 2D may loose much of its 3D information, the RGB modality.
- its geometric shape where only the 3D structure is considered, discarding all the texture information, which is the depth modality.

Based on this multimodal nature of data examples, we want to design a self-supervised task that leveraging the intermodal relations learns to abstract from the domain.

The idea behind self-supervised domain adaptation is that, along with a main classification task, is solved a self-supervised pretext task that implicitly regards intrinsic properties of the image and for which the label is not needed, hence it can be done on both source and target domain images.

As a result, the features learned after solving both tasks will be domain invariant and more suited to properly classify images from different visual domains. In the multimodal setting, exploiting the different modalities, we are able to design more robust pretext tasks that enhance intermodal relations as well.

Before delving into the implementation of a novel domain adaptation model, we repeated the results obtained in [1], trying to improve the presented results using some well known techniques in deep learning.

Our self-supervised proposed task is a task aimed at predicting

the relative difference between the zoom that has been applied independently both to RGB and depth modalities.

In our work, we show that the choice of the pretext tasks is relevant in achieving better performances of the trained model.

Another important factor is the weight initialization of the model which we found to be a key point to obtain good results. At the beginning of this project this topic was not much considered, but reading and observing the effects of Xavier initialization on our model we understood how much influence in the training it can have.

Another aspect took into account in the optimization process is the design of a proper data augmentation pipeline with the aim of reducing overfitting, improving robustness and replicating the data augmentation step involving depth images, proposed in [4].

To furtherly optimize our models we tried different networks as features extractors and a multi-task learning approach inspired by [5].

## II. RELATED WORK

One of the key references is the paper of Eitel et al. [4] which explores the possibility to use multimodal data for RGB-D object recognition. Furthermore, it proposed an intuitive data augmentation schema suited for depth images.

The leitmotiv of our work is the comparison with the model and the results proposed in [1]. In that paper the self-supervised task is implemented as the prediction of the relative rotation between RGB and depth modalities, previously independently rotated. This has inspired our work and its outstanding performances with respect to the current state of the art has encouraged us to propose an alternative that stems from it.

### A. Unsupervised Domain Adaptation

Methods for unsupervised domain adaptation in computer vision can be divided into three main classes. The first, aims to align source and target domains in some feature space. This is done by optimizing for some measurement of distributional discrepancy. One popular measurement is the maximum mean discrepancy (MMD) [6] but there exists other discrepancy-based methods as well [7, 8].

The second class contains methods based on adversarial learning. One of these, is the domain adversarial network proposed in Ganin et al. [9], in which the pretext task consists of a discriminator branch that has to classify the domain of the given image samples. However, in this case RGB-D data has not been used and the domain adaptation is instead achieved by means of a gradient reversal layer. The proposed layer backpropagates positive gradients, thus maximizing the discriminator loss that arrives to the feature extractors; this allows domain independent features to emerge. Similar works are [10], [11], [12].

The third class leverages the solution of a self-supervised tasks in parallel with the main recognition task. A self-supervised task has the benefit of being carried out without the need of

labels, therefore it can be executed on both source and target examples. This allow to learn robust domain invariant features. Although there exist examples of self-supervised domain adaptation applied to RGB-D data [13][14][15][16], they all lack something in what is our focus, namely adapting a source (synthetic) domain to a target (real) domain of RGB-D images, leveraging intermodal relations.

As far as we know, the only work addressing this issue is [1]. In this case, the pretext task consist in predicting the relative zoom of RGB and depth modalities. This enhances the multimodal relations because the pretext task learns that when there is a rotation the pixels in both modalities behave in the same way, pointing out the features that enable to recognize better across domains. In their work, they showed that these features are the shapes of the objects, rather than their texture. Eventually, in [5] is presented a novel technique of multi-task learning. In that paper it is argued that solving  $K$  pretext tasks, rather than one, can further improve the performances of the main task, entailing a better domain alignment. The proposed self-supervised task are, for instance, the prediction of the absolute rotation and the location of patch, that can be implemented as classification or as regression problems.

## III. DOMAIN ADAPTATION AND MULTI MODAL DATA

The standard Domain Adaptation problem considers paired samples  $x, y \in X \otimes Y$  sampled from different distributions, namely the source and the target distributions  $S(x, y)$  and  $T(x, y)$ , with  $S \neq T$ . The paired samples consist of items  $x$  and their labels  $y \in Y$  from the finite set  $Y = \{0, 1, \dots, L\}$ . Each item  $x_i$  is composed of multimodal data, represented with the components  $x_i^{(rgb)}$  and  $x_i^{(depth)}$ . In our setting we consider the marginal distributions  $S(x)$  and  $T(x)$  from which the training sets are constructed. In the case of the samples  $x_i \sim S(x)$ , labels are known at training time, while these are unknown for the samples  $x_i \sim T(x)$ . Therefore, the training phase consists of a supervised classification task that allows the model to extract RGB features and depth features from the source domain. In order to tackle the absence of target domain labels, a self-supervised task is employed. The self-supervision is performed jointly on source and target domain images to allow the model to learn features that are meaningful for both domains. More specifically, for each domain we consider the RGB and depth representation of each item. Then for each  $x_i^{(rgb)}$  and  $x_i^{(depth)}$  a self-supervised label is generated based on the selected type of pretext task. By training jointly on the described main and pretext tasks the aim is to use the learned domain invariant and RGB-depth features to predict labels on unknown target domain data.

## IV. DATASETS

### A. Description

For our tests we considered the SynROD and the ROD datasets. The first one was collected and created as described in [1] and consists in both RGB and depth version of artificially generated images, while ROD consists in real

RGB and depth images. These sets contain 51 classes of common objects. Out of all the categories, 47 are considered, leaving out 'lime', 'onion', 'peach', 'tomato'. Furthermore, subsets of images for each class are sampled from SynROD and ROD. Overall, following the proposed synROD split we obtain 37528 synROD train samples 7301 synROD validation images, totaling  $N_{synrod} = 44829$  samples. Regarding ROD, the considered ROD subset that is used for testing the configuration consists of  $N_{rod} = 32476$  images out of the 41877 present in the original ROD.

### B. Usage

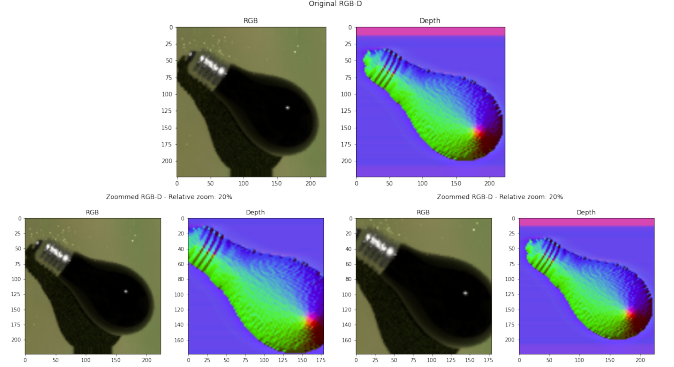
We generate the training and test sets by pairing for a given sample its RGB and depth representation along with its label. Therefore, we have a source main task training set  $S_{train} = \left\{ \left( x_i^{(rgb)} x_i^{(depth)} \right), y_i \right\}_{i=1}^{N_{synrod}}$ , a source pretext task training set  $S_{pretext} = \left\{ \left( x_i^{(rgb)} x_i^{(depth)} \right), z_i \right\}_{i=1}^{N_{synrod}}$ , a target pretext task training set  $T_{pretext} = \left\{ \left( x_i^{(rgb)} x_i^{(depth)} \right), z_i \right\}_{i=1}^{N_{rod}}$  and a target test set  $T_{test} = \left\{ \left( x_i^{(rgb)} x_i^{(depth)} \right), y_i \right\}_{i=1}^{N_{rod}}$ . Here,  $z_i$  is generated pretext task label.

## V. METHOD

### A. Pretext task

In this section we present the types of pretext task that have been tested on the dataset configuration presented in the previous paragraph. Starting from the previous work on relative rotation classification pretext task[1], we extend the experiments by proposing two different pretext task variations based on the same principles. That is, the application of relative transformations to the RGB and depth representation of the samples. We propose a relative zoom regression and an alternative classification pretext as means of learning RGB-D invariant features with the purpose of generalizing also to different domains the learned characteristics. In order to avoid loss of information caused by an high zoom applied to both images, we randomly select only one which is then zoomed in of a percentage in the range  $[0, 40]$ . Therefore, the defined label generation procedure for the pretext task consists of a first step in which one between the RGB and depth image is selected with equal probability. Then, a random number between 0 and 40 included is sampled from a uniform distribution and the selected image is zoomed-in of that amount, while the other is left as is. By following this procedure it is possible to generate a "relative zoom difference" random variable that is uniformly distributed. Since the zoom difference values follow a total ordering, the relative zoom percentage prediction can be seen as a regression problem. An alternative to this, is to predict the relative zoom difference percentage as a specific category, thus considering this setting as that of a classification problem. In our work, we consider both variations of this

**Fig. 1:** Example of 20% relative zoom



pretext problem.

### B. CNN architecture

The CNN architecture employed in this study is composed of two feature extractor networks, whose outputs are concatenated along the filter dimension and then fed into a main and a pretext branch. These are dedicated respectively to the main classification task and the given pretext task. In this analysis we propose different variations of this architecture, which consist in different feature extractor architectures combined with modifications to the pretext branch. The starting CNN architecture is the one proposed in [1] that consists of two convolutional feature extractors of which one is dedicated to RGB images and the other to the respective depth representation. For this purpose, two ResNet18 without the last global average pooling and fully connected layers have been used. The two networks are initialized with pretrained imagenet weights. Their output features are concatenated and the result is forwarded either to the main task branch or to the pretext task branch. The former will be referred to as  $M$ , while the latter as  $P$ . Regarding the structure of the  $M$  branch, it consists of a global average pooling layer, followed by a fully connected that uses batch normalization, a ReLU activation function and dropout 0.5. This branch terminates with a fully connected layer with 47 output neurons. The  $P$  branch instead consists of a sequence of two 2D convolutional layers that use batch normalization and ReLU activation function. The first of the two utilizes  $1 \times 1$  filters with stride=1, while the one that follows uses  $3 \times 3$  filters with stride=2. Padding is set to 0 for both. At the end of this sequence, a global average pooling operation is performed before feeding the result into another fully connected layer with batch normalization, ReLU and dropout 0.5. Its output ends up in the final fully connected output layer which consists of: 4 neurons with softmax output for the relative rotation task; 1 neuron for the relative zoom regression; a variable number of neurons from 1 to 41 followed by a softmax output for the relative zoom classification task. While the two resnet18 are pretrained, the  $M$  and  $P$  branches have been initialized with He normal initialization [17] in our best performing configurations. In this case, the two ReLU

activations of the pretext branch convolutional layers have been substituted with the PReLU activation proposed in [17], each setup with 1 parameter in order to limit overfitting.

The substitution of the resnets with two inceptionV3 models yields a variation of the architecture that allows us to evaluate how a wider network affects domain adaptation performances. In order to make use of the auxiliary modules that are part of the inception models, we propose two modifications that adapt them to the context of this problem. A simpler adjustment consists in replacing their fully connected output layer with one dedicated to the main task and another to the pretext. This structure mimics exactly that of the final output layer of the  $M$  and  $P$  branches. A more complex modification instead consists in removing the global average pooling and fully connected output layer from the auxiliary modules to allow for the concatenation of their output filters. These convolutional outputs are then fed into a fc-softmax output layer dedicated to the main task and into a fc-output layer reserved for the pretext task, thus mimicking the structure of the  $M$  and  $P$  branches. The main and pretext branches attached to the concatenation of the output of the RGB-inception and of the Depth-inception maintain the same structure proposed in the previous multi resnet18 architecture.

### C. Initialization

1) *Default Initialization:*

2) *Xavier Initialization:*

3) *Kaiming Initialization:* After the remarkable initialization method proposed by [18], that was not designed for generic activation function, a further step was made by [17] in which were introduced a generalization or rectified activation functions (ReLU), called Parametric Rectified Linear Unit (PReLU) and a novel initialization schema, originally thought for convolutional layers, tailored for rectified units.

### D. Data Augmentation

A common practice in deep learning to increase the capability of a model to generalize and to increase its robustness is indeed data augmentation [19].

Our proposed augmentation technique is tailored separately for each modality, trying to leverage their peculiar properties. The augmentation carried out on RGB images is based on the assumption that in a couple RGB-D the depth image is immune to changes in lighting conditions, such as the change of sun illumination during the day. For this reason, exclusively on the RGB image it has been applied color jittering, remembering that the depth is just the result of a colorization preprocessing, and its colors does not have the same semantic meaning.

As stated in [4], depth cameras are not perfect and suffer from light reflections, thus real depth images are subject to noise and occlusions phenomena. These are a missing information comparable to random erasure of some areas. Therefore, in the cited paper is proposed an advanced augmentation procedure applied to synthetic depth images, which are too "clean",

consisting in erasing some pixels to simulate real disturbs.

We decided to adopt a similar augmentation strategy, implementing the Cutout (or random erasing) proposed in [20], encouraged by its outstanding results. The cutout is performed with a 50% probability, a variable scale between (0.02, 0.08) and a variable aspect ratio between (0.03, 0.33)

Eventually, for each original RGB-D image another is created by flipping it horizontally. The color jittering augmentation is performed only on one of the two.

### E. Multiple pretexts

As presented in [5], using more than one self-supervised task may induce a stronger alignment among source and target domains. This is because the solution of each pretext tasks involves learning some specific domain invariant features. Having more pretext tasks may produce even more diverse and equally robust features. Following their idea and inspired by their results, we tried to implement a model having two pretext tasks:

- prediction of the relative rotation between RGB and depth, as presented in [1].
- prediction of the relative zoom, our proposed variation.

The zoom prediction task is implemented as a classification and as a regression problem, in different trials. Our expectations are that using this configuration will leverage the cooperation between both auxiliary tasks, leading to better results than the ones obtained implementing just one or another.

### F. Training

The proposed training procedure requires the computation of the predictions for the main classification task in the context of the source domain and of the pretext regression task in the source and target domains. Therefore, the losses that are employed in order to train the model are the mean squared error and the cross-entropy loss. Due to the fact that no target domain label are available at training time, an entropy minimization regularization procedure is also used for the purpose of learning partial information also on the target domain main classification task.

Let  $\theta_M$ ,  $\theta_P$  and  $\theta_F$  be the parameters of the main branch, the pretext branch and of the feature extractors respectively. Let

$$\hat{y}_i^s = M(F(\hat{S}); \theta_M, \theta_F, S)$$

be the output of the main branch, received after the rgb and depth images of the source batch  $\hat{S}$  have been fed into the respective feature extractors and then concatenated. Let

$$\hat{z}_i^s = P(F(\hat{S}_{rot}); \theta_P, \theta_F, \hat{S}_{rot})$$

be the output of the pretext branch, received after the rgb and depth images of the source rotated batch  $\hat{S}_{rot}$  have been fed into the respective feature extractors and then concatenated. After defining the outputs related to the target domain with only a change in the notation of the previous from  $s$  to  $t$ , we



denote the losses computed by means of the defined values as follows:

$$L_{main} = -\frac{1}{N_s} \sum_{i=1}^{N_s} y_i^{ss} \log(\hat{y}_i^s)$$

$$L_{reg\_pretext} = \frac{1}{N_s} \sum_{i=1}^{N_s} (z_i^s - \hat{z}_i^s)^2 + \frac{1}{N_t} \sum_{i=1}^{N_t} (z_i^t - \hat{z}_i^t)^2$$

$$L_{clf\_pretext} = -\frac{1}{N_s} \sum_{i=1}^{N_s} z_i^s \log(\hat{z}_i^s) - \frac{1}{N_t} \sum_{i=1}^{N_t} z_i^t \log(\hat{z}_i^t)$$

The complete loss consists of the sum of the main task loss and the two pretext tasks losses weighted by a factor  $\lambda$ :  $L_{tot} = L_{main} + \lambda L_{pretext}$ . As mentioned before, regularization techniques such as entropy minimization and weight decay are also employed in order to reduce overfitting. Both are implemented and weighted in the parameter update formula. The entropy minimization is weighted with a coefficient  $\tau$  set to 0.1, while the weight decay  $wd$  value is set to  $5e-2$ . While these two hyperparameters are kept constant for each tested configuration, the learning rate and  $\lambda$  have been fine tuned by means of grid search.

---

**Algorithm** RGB-D neural network training

---

```

1: Input
2:   source set  $S = \{(x_i^{(rgb)}, x_i^{(depth)}), y_i\}_{i=1}^{N_{synrod}}$ 
3:   target set  $T = \{(x_i^{(rgb)}, x_i^{(depth)})\}_{i=1}^{N_{rod}}$ 
4:   network
5: Output
6:   test set predictions  $\hat{y}$ 
7: procedure TRAINING_EPOCH( $network, S, T$ )
8:   for  $i = 1$  to  $N_{batches}$  do
9:     load main task source batch  $\hat{S}$  from  $S$ 
10:    load target batch  $\hat{T}$  from  $T$ 
11:    load pretext batch  $\hat{S}_{rot}$  from  $S$ 
12:    load pretext batch  $\hat{S}_{rot}$  from  $T$ 
13:
14:    do network forward pass
15:
16:    compute main source loss  $L_{main}$ 
17:    compute source pretext loss  $L_{sp}$ 
18:    compute target pretext loss  $L_{tp}$ 
19:    compute entropy loss  $L_{ent}$ 
20:     $Loss \leftarrow L_{main} + \lambda(L_{sp} + L_{tp})$ 
21:
22:    update all parameters  $\theta$  of the network
23:     $\theta \leftarrow \theta - \eta \nabla_{\theta}(Loss(\theta) + \tau L_{ent}(\theta)) - wd \cdot \theta$ 
24: procedure TEST( $network, T$ )
25:   for  $i = 1$  to  $N_{test\_batches}$  do
26:     load test batch from  $T$ 
27:     predict target main task labels  $\hat{y}$ 

```

---

The empirical loss computed in the presented training procedure is composed of three contributions. The main loss  $L_{main}$ , calculated by forwarding a source batch through the two feature extractors and the  $M$  branch of the network. Therefore, this loss does not depend on the  $P$  branch parameters. The  $L_{sp}$  and  $L_{tp}$  losses are computed by forwarding the source or target pretext batch into the two feature extractors and then sending their concatenated output to the  $P$  branch. Thus, these two contributions do not depend on  $M$  branch parameters. The vector  $\theta$  contains all parameters of the network.

## VI. EXPERIMENTS

In this section we present our experiments during training. It is worth of notice that we tried two ways of pre-processing the images: the first one normalizing them using ImageNet mean and standard deviation for both RGB and depth modalities, while the second using ImageNet mean and standard deviation for RGB images and a custom mean and standard deviation computed on the depth dataset for the depth modality. We decided to try this because RGB and depth images belong to different domains and their shift is not slight. We calculated mean and standard deviation of the depth dataset and found that they are very different from ImageNet. In this section we used both ImageNet only normalization and depth custom normalization. Best results are available in table 1.

### A. Baselines

The performance improvement obtained by means of domain adaptation and multi modal data on the target domain main classification task has been evaluated by considering various baseline models that do not make use of a pretext task. We used five types of baselines: RGB only, depth only, RGB-D, RGB-D e2e and RGB-D e2e (main head). All these baseline experiments were performed using the complete datasets with all the 51 classes. Another consideration about these baselines is that the layers added were all initialized using Xavier uniform initialization [18] which during our work we found to be not the most efficient on fully connected layers. All the baselines were implemented following Robbiano et al. [1] paper without using any dataset splits or less than 51 classes.

1) *RGB only*: The network architecture used to evaluate this baseline is a single pretrained ResNet18 with the last fully connected with 51 output neurons. During training only RGB images from the source dataset (synROD) were feed to the network and at the end the model was tested on the target RGB dataset (ROD). This is a standard non Domain Adaptation algorithm which however didn't cut a poor figure because it was the best baseline model. For hyperparameters tuning we fine tuned learning rate in [0.0001, 0.001, 0.01], batch size [32, 64] and weight decay [0.05, 0.005] using 20 epochs and step size 15. After tuning we also re-trained the best model found 5 times to see the mean and the consistency of the training.

2) *Depth only*: For this part we used the same architecture as for RGB only but this time we also tried training the model normalizing the depth images using ImageNet means or our custom depth means. The results are very similar in accuracy, but ImageNet gave better results. As expected, training a model on depth only features didn't lead to good results. That is because depth images are not much detailed and the feature extractors find it very difficult to obtain domain invariant features using only depth representation of the objects. As before we performed a grid search algorithm using the same hyperparameters as for the RGB only section and then re-trained the best models found 5 times each.

3) *RGB-D*: The network used for the RGB-D baseline is the combination of the RGB only and depth only networks where the last fully connected layers are removed and the concatenation of the two extracted features is used as input for a new fully connected layer. For this method the two feature extractors were already trained and frozen and only the last fully connected was trained. We used our best models found in the previous RGB only and depth only as the two feature extractor. We optimized hyperparameters also in this method but with batch size fixed at 64.

4) *RGB-D e2e and RGB-D e2e (main head)*: In RGB-D e2e the network is the same as RGB-D but the feature extractors are not frozen while in RGB-D e2e (main head) baseline the network is the same as RGB-D, but instead of the last fully connected, the extracted and concatenated RGB and depth features are fed to the main branch of the network. The feature extractors are not frozen and start the training pretrained on ImageNet. In this way the model is trained in an end to end fashion. Also for this method we ran a grid search for best hyperparameters.

## B. Variation

The proposed architecture uses its main task alongside a pretext task as a means of achieving better domain adaptation performance and generalization capability of the learned domain independent features. In our experiments we tested various configurations and tuned the learning rate and lambda hyperparameters in order to achieve the best performing configuration. The model has been trained on 20 epochs, with batch size of 64. Regarding the optimizer, SGD with momentum 0.9 has been employed. The hyperparameters considered in the analysis are learning rates in the range between  $[1e-4, 1e-2]$  and values of  $\lambda \in [0.001, 1.0]$ .

## VII. CONCLUSIONS

### REFERENCES

- [1] M. Reza Loghmani, L. Robbiano, M. Planamente, K. Park, B. Caputo, and M. Vincze, "Unsupervised domain adaptation through inter-modal rotation for rgb-d object recognition," 2020.
- [2] K. Saenko, B. Kulis, M. Fritz, and T. Darrell, "Adapting visual category models to new domains," 2010.
- [3] F. M. Carlucci, P. Russo, and B. Caputo, "(de)2co: Deep depth colorization," 2018.
- [4] A. Eitel, J. T. Springenberg, L. Spinello, M. Riedmiller, and W. Burgard, "Multimodal deep learning for robust rgb-d object recognition,"

Baselines		
<b>RGB</b>	imagenet	50.93
<b>Depth</b>	imagenet	12.50
	depth custom	11.71
<b>RGB-D</b>	imagenet	45.90
	depth custom	50.00
<b>RGB-De2e</b>	imagenet	44.10
	depth custom	43.16
<b>RGB-De2e (main head)</b>	imagenet	46.00
	depth custom	42.21

TABLE I: Best baseline results

- [5] Y. Sun, E. Tzeng, T. Darrell, and A. A. Efros, "Unsupervised domain adaptation through self-supervision," 2019.
- [6] M. Long, Y. Cao, J. Wang, and M. I. Jordan, "Learning transferable features with deep adaptation networks," 2015.
- [7] B. Sun, J. Feng, and K. Saenko, "Return of frustratingly easy domain adaptation," 2016.
- [8] R. Xu, G. Li, and J. Yang, "Larger norm more transferable: An adaptive feature norm approach for unsupervised domain adaptation," 2019.
- [9] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," 2016.
- [10] P. Russo, M. Carlucci, T. Tommasi, and B. Caputo, "From source to target and back: symmetric bi-directional adaptive gan," 2018.
- [11] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," 2017.
- [12] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," 2014.
- [13] L. Spinello and K. O. Arras, "Leveraging rgb-d data: Adaptive fusion and domain adaptation for object detection," in *ICRA. IEEE*, p. pp. 4469–4474, 2012.
- [14] J. Hoffman, S. Gupta, J. Leong, S. Guadarrama, and T. Darrell, "Crossmodal adaptation for rgb-d detection," in *ICRA. IEEE*, p. pp. 5032–5039, 2016.
- [15] X. Li, J. Fang, J. Zhang, J. Wu, and T. Darrell, "Domain adaptation from rgb-d to rgb images," *Signal Processing*, vol. 131, p. pp.27–35, 2017.
- [16] W. Jing and Z. Kuang, "Unsupervised domain adaptation learning algorithm for rgb-d staircase recognition," *URL http://arxiv.org/abs/1903.01212*, 2019.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," 2015.
- [18] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," 2010.
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," 2012.
- [20] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," *URL http://arxiv.org/abs/1708.04896*, 2017.